UVA CS 4774: Machine Learning

Part 2: Course Session QA

Related:

- Part 1 [HERE]
- Course Announcement Page [HERE]

Dr. Yanjun Qi

University of Virginia

Department Of Computer Science

09/30

09/30/2025 Assignments

- HW2 is due this coming Sunday midnight!
 - Using HW1 code pieces as components;
 - If you struggle with HW1, please contact TA @Haochen ASAP
- HW1 grading is work-in-progress,
 - Grades will be released by next Tuesday class time
 - We posted the guide from TA in Canvas
- Course vote:
 - New Survey that needs your vote on
 - 1. back to lecture in-person twice a week?
 - 2. If not, best way to use the in-person session:
 - Quiz to continue
 - + Project discussions
 - Interested in Shark Tank alike setup? idea screening, pitch talk, demo ...

Project Process

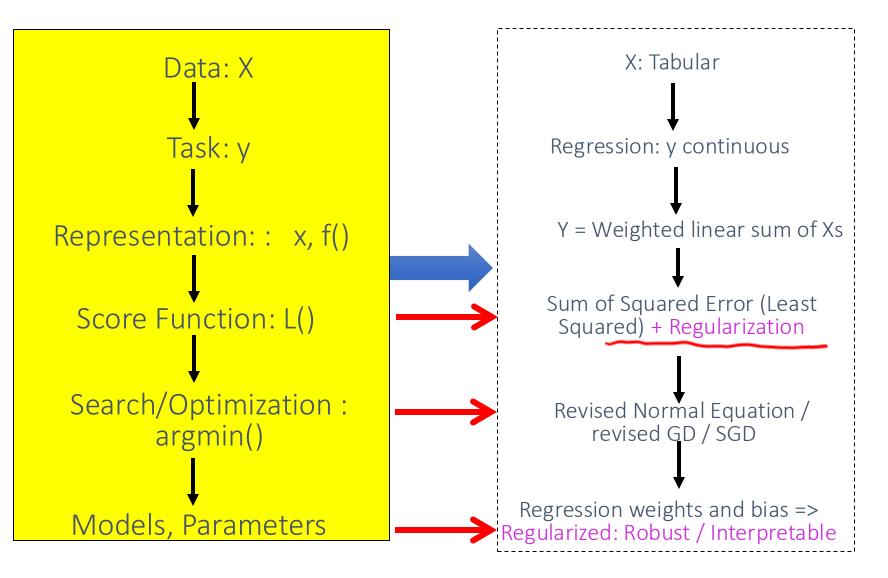


- Format:
 - Team, individual ...
 - Shark Tank alike Screening? –
 https://en.wikipedia.org/wiki/Shark_Tank
 - Next week Idea collection
- Final deliverables:
 - (1) Code (Github PR to course project repo)
 - (2) Poster presentation class wide (Date: TBD)
 - (3) Video Demo (TBD)

09/30/2025 Roadmap

- •TA to go over HW1
- One UVA ML club to introduce their setups and projects
- •Q5
- Review Q4
- Review QA for L5-L7

L7: Regularized multivariate linear regression



10/22/2025

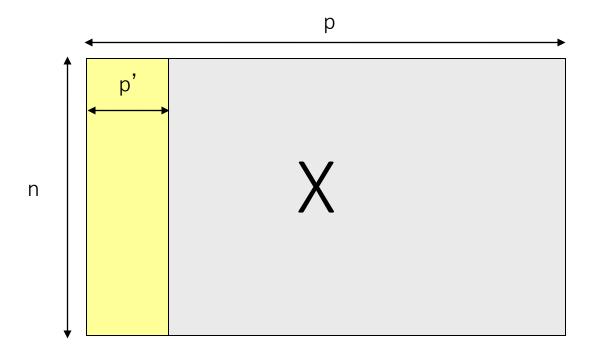
L7: Regularized multivariate linear regression

We aim to make our trained model

•1. Generalize Well

- 2. Computational Scalable and Efficient
- 3. Trustworthy: Robust / Interpretable
 - Especially for some domains, this is about trust!

Large p, small n: How? $\uparrow \rightarrow \uparrow' \Rightarrow \downarrow easy to understand$



Regularized multivariate linear regression

• Model:
$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$$

• LR estimation:

$$\underset{?}{\operatorname{arg\,min}} \sum \left(Y - \hat{Y} \right)^{2}$$

• LASSO estimation:

$$\underset{i=1}{\operatorname{arg\,min}} \sum_{i=1}^{n} \left(Y - Y \right)^{2} + \lambda \sum_{j=1}^{p} \left| \beta_{j} \right|$$

• Ridge regression estimation:

$$\underset{i=1}{\operatorname{arg\,min}} \sum_{i=1}^{n} \left(Y - Y \right)^{2} + \lambda \sum_{j=1}^{p} \beta_{j}^{2}$$

$$?$$

Error on data

Regularization

Ridge Regression / L2 Regularized Regression

$$\boldsymbol{\beta}^* = \left(\boldsymbol{X}^T \boldsymbol{X} \right)^{-1} \boldsymbol{X}^T \boldsymbol{\bar{y}}$$

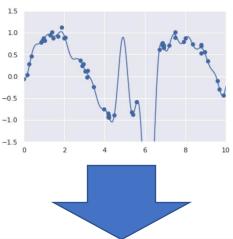


• If not invertible, a classical solution is to add a small positive element to diagonal

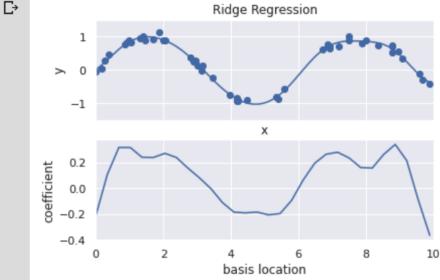
$$\boldsymbol{\beta}^* = \left(\boldsymbol{X}^T \boldsymbol{X} + \lambda \boldsymbol{I} \right)^{-1} \boldsymbol{X}^T \, \boldsymbol{\bar{y}}$$

Overfitting: Can be Handled by Regularization

A regularizer is an additional criteria to the loss function to make sure that we don't overfit. It's called a regularizer since it tries to keep the parameters more normal/regular







WHY and How to Select λ ?

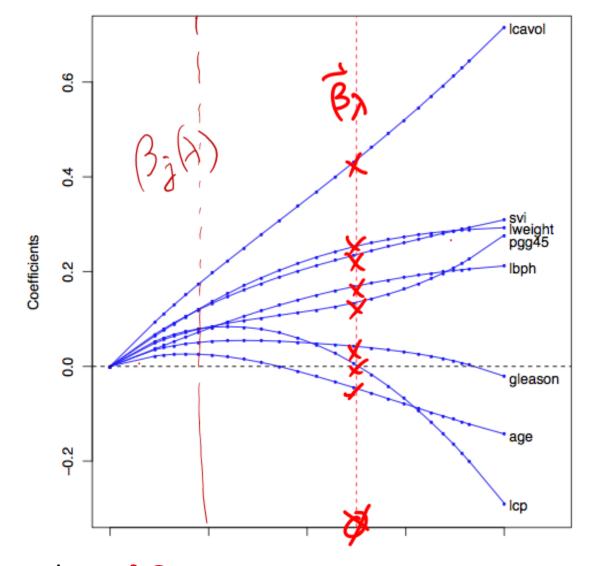
- 1. Generalization ability
 - → k-folds CV to decide
- 2. Control the bias and Variance of the model (details in future lectures)

L2: Squared weights penalizes large values more

L1: Sum of weights will penalize small values more

$$egin{array}{c} eta_j^2 \ eta_j \end{array}$$

Regularization path of a Ridge Regression



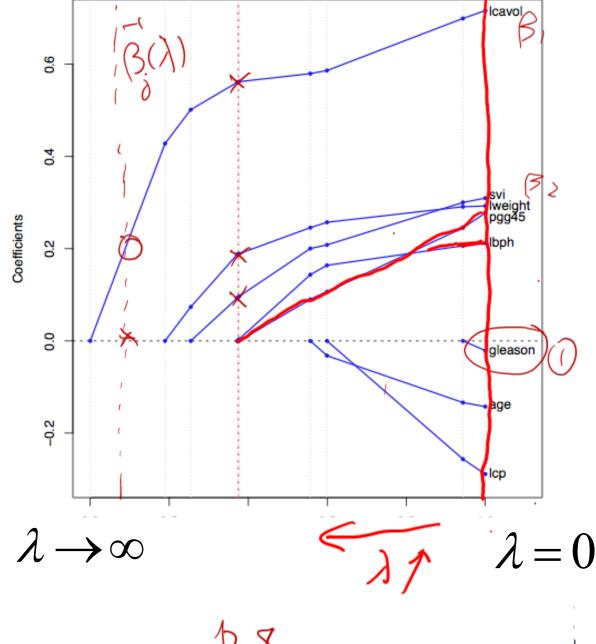
WHY and How to Select λ ?

$$\lambda \to \infty$$

$$\lambda = 0$$

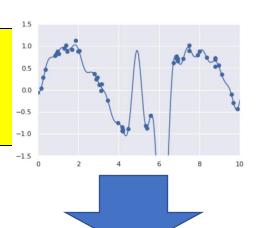
Regularization path of a Lasso Regression

when varying λ , how β_i varies.



Overfitting: Can be Handled by Regularization

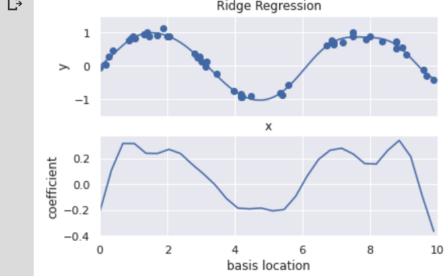
A regularizer is an additional criteria to the loss function to make sure that we don't overfit. It's called a regularizer since it tries to keep the parameters more normal/regular



code-run:

https://github.com/qiyanju n/2025Fall-UVA-CS-MachineLearningDeep/blo b/main/notebook/L7 regul arizedRegression 06 Linea r Regression.ipynb





(Extra) Lasso (least absolute shrinkage and selection operator) / Squared Loss+L1

- The lasso is a shrinkage method like ridge, but acts in a nonlinear manner on the outcome y.
- The lasso is defined by

$$\hat{\beta}^{lasso} = \operatorname{argmin}(y - X \beta)^{T} (y - X \beta)$$

$$\operatorname{subject} \beta_{j}^{lasso} = s$$

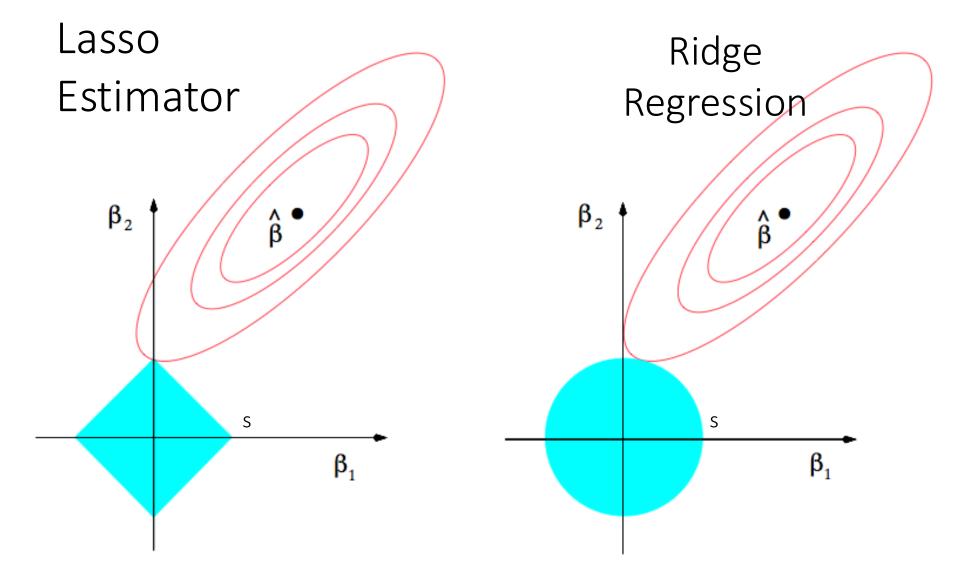
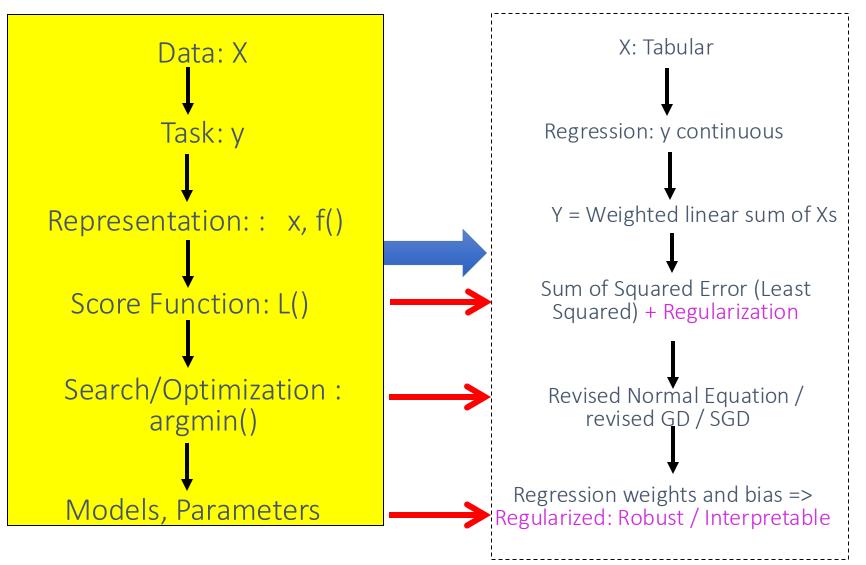


FIGURE 3.11. Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions $|\beta_1| + |\beta_2| \le t$ and $\beta_1^2 + \beta_2^2 \le t^2$, respectively, while the red ellipses are the contours of the least squares error function.

$$\min J(\beta) = \sum_{i=1}^{n} \left(Y - \hat{Y} \right)^{2} + \lambda \left(\sum_{j=1}^{p} \beta_{j}^{q} \right)^{1/q}$$

Today: Regularized multivariate linear regression



10/22/2025

More: A family of shrinkage estimators

$$\beta = \operatorname{arg\,min}_{\beta} \left[\left| y_{i} - x_{i}^{T} \beta \right|^{2} \right]$$
subject $\beta = \left| \beta_{j} \right|^{q} \le s$

• for q >=0, contours of constant value of $\sum_j \left| \beta_j \right|^q$ are shown for the case of two inputs.

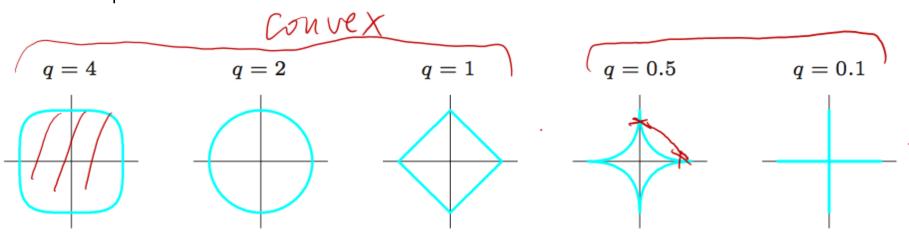
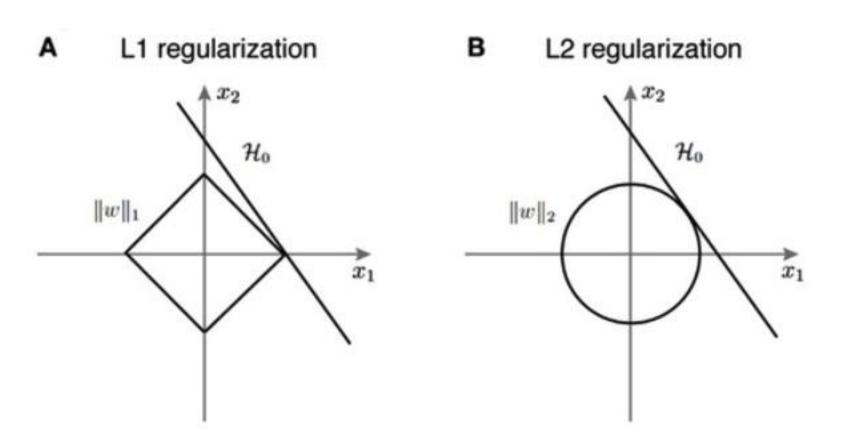


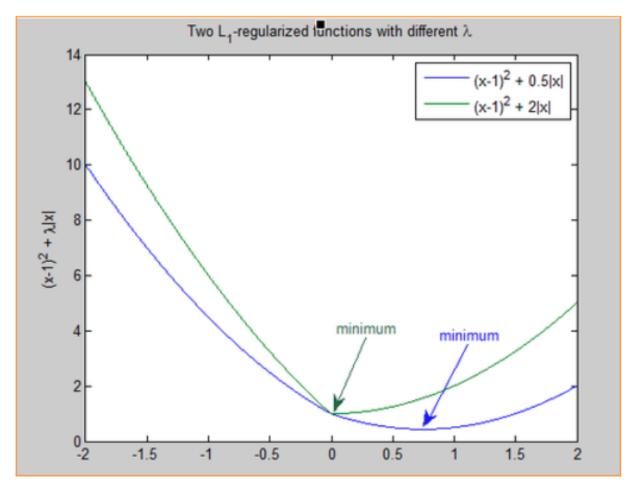
FIGURE 3.12. Contours of constant value of $\sum_{j} |\beta_{j}|^{q}$ for given values of q.



due to the nature of L_1 norm, the viable solutions are limited to corners, which are on a few axis only - in the above case x1. Value of x2 = 0. This means that the solution has eliminated the role of x2, leading to sparsity

 L_1 -regularized loss function $F(x) = f(x) + \lambda ||x||_1$ is non-smooth. It's not differentiable at o. Optimization theory says that the optimum of a function is either the point with o-derivative or one of the irregularities (corners, kinks, etc.). So, it's possible that the optimal point of F is o even if o isn't the stationary point of f. In fact, it would be o if λ is large enough (stronger regularization effect). Below is a graphical illustration.

http://www.quora.com/What-is-the-difference-between-L1-and-L2-regularization



In mathematics, particularly in calculus, a stationary point or critical point of a differentiable function of one variable is a point of the domain of the function where the derivative is zero (equivalently, the slope of the graph at that point is zero).

Coordinate descent based Learning of Lasso

Coordinate descent
(WIKI)→ one does
line search along one
coordinate direction
at the current point in
each iteration.

One uses different coordinate directions cyclically throughout the procedure.

1. Initialize
$$\beta$$

2. Repeat until converged

3. For $j = 1, 2, ..., P$ do

$$a_{j} = 2 \sum_{i=1}^{m} x_{ij}^{2}$$

$$c_{j} = 2 \sum_{i=1}^{m} x_{ij} (y_{i} - x_{i}^{T}\beta + x_{ij}\beta_{j})$$

$$if c_{j} < -\lambda$$

$$\beta_{j} = (c_{j} + \lambda)/\alpha_{j}$$

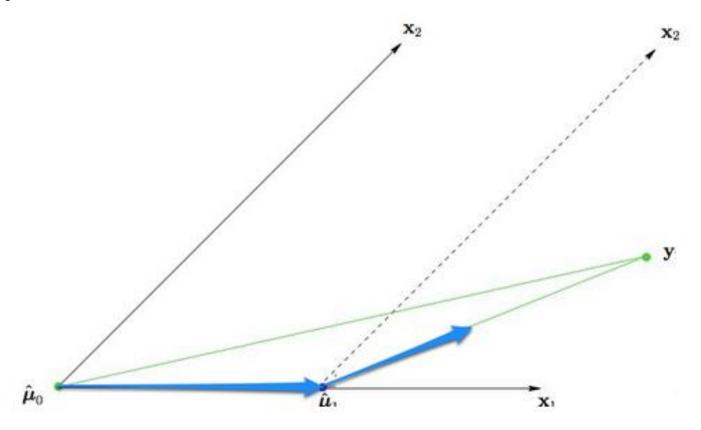
$$else if, c_{j} > \lambda$$

$$\beta_{j} = (c_{j} - \lambda)/\alpha_{j}$$

$$else soft-thresholding$$

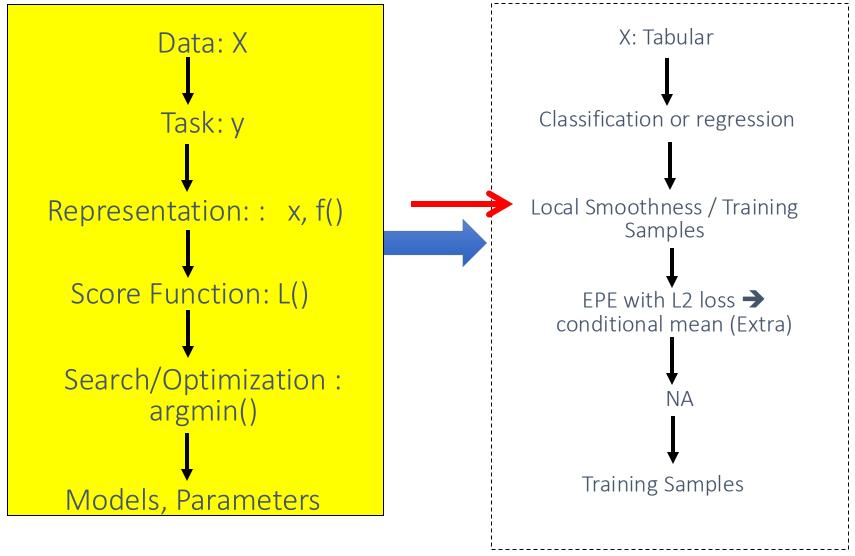
$$\beta_{j} = 0$$

Least Angle Regression (LARS) (State-of-the-art LASSO solver)



http://statweb.stanford.edu/~tibs/ftp/lars.pdf

L8: K-nearest-neighbor(regressor or classifier)



Yanjun Qi @ UVA CS

Code run: https://github.com/qiyanjun/2025Fall-UVA-CS- MachineLearningDeep/blob/main/notebook/L8 Knearest.ipynb

```
[42] # import regressor
    from sklearn.neighbors import KNeighborsRegressor
    # instantiate with K=5
    knn = KNeighborsRegressor(n_neighbors=5)
    # fit with data
    knn.fit(X, y)
```

```
###
Xfit = np.linspace(3, 10, 1000).reshape(-1, 1)
yfit =knn.predict(Xfit)

# Plot outputs
plt.scatter(X, y, color='red')
plt.plot(Xfit, yfit, color='blue', linewidth=3)

plt.xticks(())
plt.yticks(())
plt.show()
```

Ľ⇒

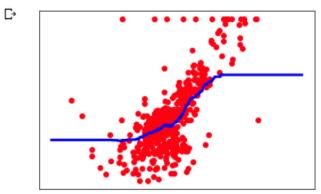
```
[44] # import regressor
    from sklearn.neighbors import KNeighborsRegressor
    # instantiate with K=5
    knn = KNeighborsRegressor(n_neighbors=100)
    # fit with data
    knn.fit(X, y)
```

KNeighborsRegressor(algorithm='auto', leaf_size=3(metric_params=None, n_jobs=Nor weights='uniform')

```
###
Xfit = np.linspace(3, 10, 1000).reshape(-1, 1)
yfit =knn.predict(Xfit)

# Plot outputs
plt.scatter(X, y, color='red')
plt.plot(Xfit, yfit, color='blue', linewidth=3)

plt.xticks(())
plt.yticks(())
plt.show()
```



K Nearest neighbor (Testing Mode)

It Needs:

- The set of stored training samples
- Distance metric to compute distance between samples
- 3. The value of k, i.e., the number of nearest neighbors to retrieve

Training Mode:

 (Naïve) version: DO NOTHING !!!!

Testing Model: To classify unknown sample:

- Step1: Compute distance to all training records
- Step2: Identify k nearest neighbors
- Step3: Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

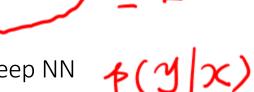
We can divide the large variety of supervised classifiers into roughly three major types.

1. Discriminative

directly estimate a decision rule/boundary

e.g., support vector machine, decision tree,

e.g. logistic regression, neural networks (NN), deep NN



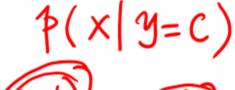
2. Generative:

build a generative statistical model

e.g., Bayesian networks, Naïve Bayes classifier



- Use observation directly (no models)
- e.g. K nearest neighbors







More Less complex complex polynomial Regression Yanjun Qi @ UVA CS

Model Selection for Nearest neighbor classification

- Choosing the value of k:
 - If k is too small, sensitive to noise points
 - If k is too large, neighborhood may include points from other classes
- •Bias and variance tradeoff

 •A small neighborhood → large variance → unreliable estimation
 •A large neighborhood → large bias → inaccurate estimation under its

We aim to make our trained model

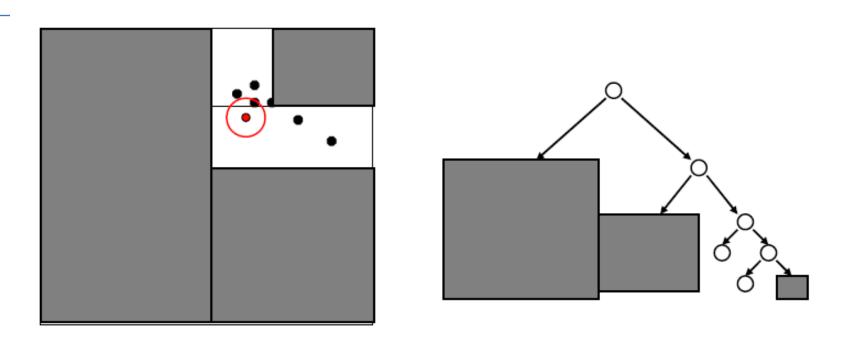
•1. Generalize Well

- 2. Computational Scalable and Efficient
- 3. Trustworthy: Robust / Interpretable
 - Especially for some domains, this is about trust!

Computational Time Cost

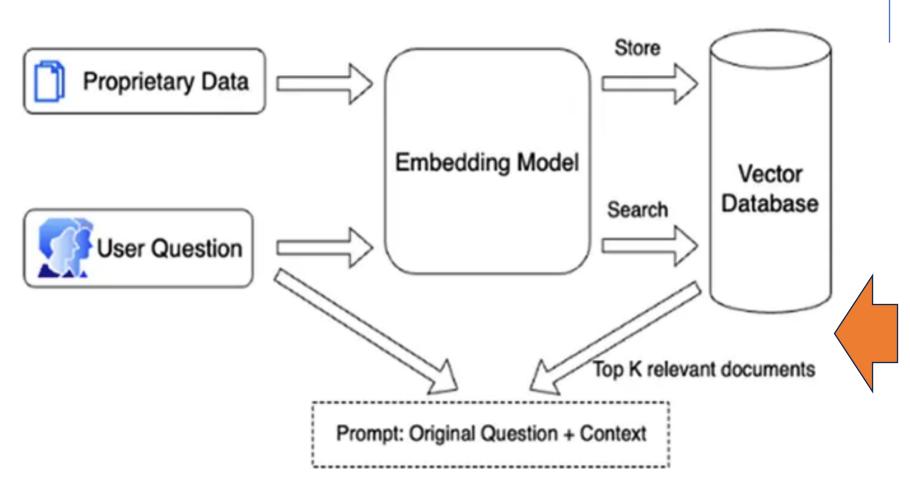
	Train (n)	Test (m=1)
Linear Regte Solon	$O(np^2+p^3)$	$O(p)^{\widehat{\hat{y}} = \beta^T x}$
KNN Reg	(I) P	0 (np)+ 0 (sort n-k) 20,000

NN Search by KD Tree



Using the distance bounds and the bounds of the data below each node, we can prune parts of the tree that could NOT include the nearest neighbor.

KNN as the Most critical component in Retrieval Augmented Generation System, e.g.:



09/30/2025 Roadmap

- •TA to go over HW1
- One UVA ML club to introduce their setups and projects
- •Q5
- Review Q4
- Review QA for L5-L7

10/07

10/07 /2025 Assignments

- HW2 is grading started ..
 - → HW2 key walk-through next Thursday online zoom
- HW3 will get posted by tomorrow
 - Deep NN on Imaging task / Kera / mostly about learning modern DNN library
 - Programming + QA (like calculating marginal prob...)
- Next Tuesday is reading day
 - → we will host makeup-Quiz Q7 next Thursday online

- Course format survey:
 - https://forms.gle/PkWGMkwHhawqf8QR8
 - Now go over the results:

Project Process



• Format:

- Team (1~4 students)
- Shark Tank alike Screening https://en.wikipedia.org/wiki/Shark_Tank
- This week: signup sheet for your team's screening sessions!
- Next week: Initial project idea collecting!
- TA Guangzhi will announce the process and signup sheet URL!

Final deliverables:

- (1) Code (Github PR to course project repo)
- (2) Poster presentation class wide (Date: 12/09 TBD)
- (3) Video Demo (after final exam)

10/07 /2025 Roadmp



Review L5-L8 questions



Quick Review L9-L10



Review Q5



Then Q6

quite disturbing for staying students around the right after quiz period

So we will host quiz after review / before project screening for all coming in-person sessions

Questions on L5-L8

Set 1: Bias-Variance, Overfitting, and Model Complexity

- •How do bias and variance contribute to generalization error, and how do we find the "sweet spot" without knowing the true distribution?
- •What are practical indicators of underfitting vs. overfitting (from graphs, learning curves, or error plots), and how do we fix each?
- •How does cross-validation (choice of K) approximate generalization error, and what are the trade-offs (bias vs variance, LOOCV vs k-fold)?
- •Why does zero training error often generalize poorly, and how is this linked to variance?
- •Would we ever prefer high bias or high variance, and how do we reduce one while controlling the other?

Set 2: Regularization (LASSO, Ridge, Elastic Net, Generalizations)

- •What are the key differences between L1 (LASSO), L2 (Ridge), and Elastic Net in terms of sparsity, robustness, computational cost, and when to use each?
- •Why do L1 penalties set coefficients to zero, while L2 does not? What happens when p>np>n or when features are highly correlated?
- •How does the choice of λ affect bias–variance, and how do we select it (cross-validation, validation curves)?
- •Are there equivalent closed-form solutions for LASSO like Ridge has? Why are L1 and L2 chosen—what about higher-order penalties?

39

•When is Elastic Net preferable (e.g., grouped correlated features), and can we always default to it?

Dr. Yanjun Qi / UVA CS

Questions on L5-L8

Set 3: k-Nearest Neighbors (kNN) and Instance-Based Learning

- •How do we pick the best k (odd vs even, weighted vs unweighted, trade-offs with noisy data)?
- •What is the computational cost of kNN (sorting term, memory cost), and can it overfit?
- •How does the distance metric affect performance, and how are ties handled in classification?
- •What are the advantages/disadvantages of kNN vs gradient descent or regularized linear models?
- •In practice, how large must the dataset be to offset outliers, and is kNN more effective for regression or classification?

Set 4: Maximum Likelihood Estimation (MLE) and Probability Foundations

- •Why do we usually maximize the log-likelihood instead of the likelihood itself, and how does this connect to squared error in linear regression?
- •How does MLE extend from discrete distributions (e.g., coin flips) to continuous (e.g., Gaussians)?
- •Why is the MLE for Bernoulli just the sample proportion, and what happens with small samples or noisy data?
- •What makes MLE consistent and efficient, and are there situations where maximum likelihood may not yield the most "ideal" parameter?
- •How does the bias–variance decomposition change with different loss functions (e.g., 0–1 loss, Laplace errors)?

10/07 /2025 Roadmp



Review L5-L8 questions



Quick Review L9-L10



Review Q5



Then Q6

quite disturbing for staying students around the right after quiz period

So we will host quiz after review / before project screening for all coming in-person sessions

Lecture 10: Maximum Likelihood Estimation (MLE)

- Probability Review
 - The big picture
 - Events and Event spaces
 - Random variables
 - Joint probability, Marginalization, conditioning, chain rule, Bayes Rule, law of total probability, etc.
 - Structural properties, e.g., Independence, conditional independence
 - Maximum Likelihood Estimation

If hard to directly estimate from data, most likely we can estimate

- 1. Joint probability
 - Use Chain Rule

$$\phi(A,B) = \phi(B) \phi(A|B)$$

- 2. Marginal probability
 - Use the total law of probability
- 3. Conditional probability
 - Use the Bayes Rule

$$P(B) = P(B, A) + P(B, A)$$
 $P(B, A) \neq P(B, A)$
 $P(B, A \cup A) \neq P(B, A)$

$$P(A|B)$$

 $P(B|A) = \frac{P(A|B)P(B)}{P(A)} = \frac{P(A|B)P(B)}{P(A)}$

One Example

Assume we have a dark box with 3 red balls and 1 blue ball. That is, we have the **set** {**r**,**r**,**r**,**b**}. What is the probability of drawing 2 red balls in the first 2 tries?

$$P(B_{1}=r,B_{2}=r) = P(B_{1}=r) P(B_{2}=r|B_{1}=r) = \frac{1}{2}$$

$$P(B_{2}=r) = P(B_{1}=r,B_{2}=r) + P(B_{1}=b,B_{2}=r)$$

$$P(B_{1}=r|B_{2}=r) = P(B_{1}=r,B_{2}=r)$$

$$P(B_{2}=r|B_{2}=r) = P(B_{1}=r,B_{2}=r)$$

MLE idea is to

- \checkmark assume a particular model with unknown parameters, θ
- ✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(Z_i|\theta)$
- ✓ We have observed a set of outcomes in the real world.
- ✓ It is then possible to choose a set of parameters which are most likely to have produced the observed results.

$$\hat{\theta} = \underset{\theta}{argmax} P(Z_1...Z_n|\theta)$$
 Likelihood

This is maximum likelihood.

In most cases this scorer is both consistent and efficient.

$$log(L(\theta)) = \sum_{i=1}^{n} log(P(Z_i|\theta))$$
 Log-Likelihood

It is often convenient to work with the Log of the likelihood function.

Deriving the Maximum Likelihood Estimate for Bernoulli

$$\begin{aligned} &\log(L(p)) \\ &= \log\left[\prod_{i=1}^{n} p^{z_i} (1-p)^{1-z_i}\right] \\ &= \sum_{i=1}^{n} (z_i \log p + (1-z_i) \log(1-p)) \\ &= \log p \sum_{i=1}^{n} z_i + \log(1-p) \sum_{i=1}^{n} (1-z_i) \\ &= \log p + (n-x) \log(1-p) \end{aligned}$$

Observed data → x heads-up from n trials

Deriving the Maximum Likelihood Estimate for Bernoulli

$$\frac{1}{p} = \frac{1}{p} \left(-x \log(p) - (n-x) \log(1-p)\right)$$

$$\frac{dl(p)}{dp} = -\frac{x}{p} - \frac{-(n-x)}{1-p} = 0$$

$$0 = -\frac{x}{p} + \frac{n-x}{1-p}$$

$$0 = \frac{-x(1-p)+p(n-x)}{p(1-p)}$$

$$Q = -x + px + pn - px$$

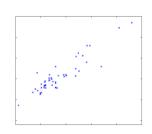
$$\Theta = -x + pn$$

Minimize the negative log-likelihood

→ MLE parameter estimation

$$\hat{p} = \frac{x}{n}$$
 i.e. Relative frequency of a binary event

DETOUR: Probabilistic Interpretation of Linear Regression



 Let us assume that the target variable and the inputs are related by the equation: $RV \in N(0, 0^2)$

$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

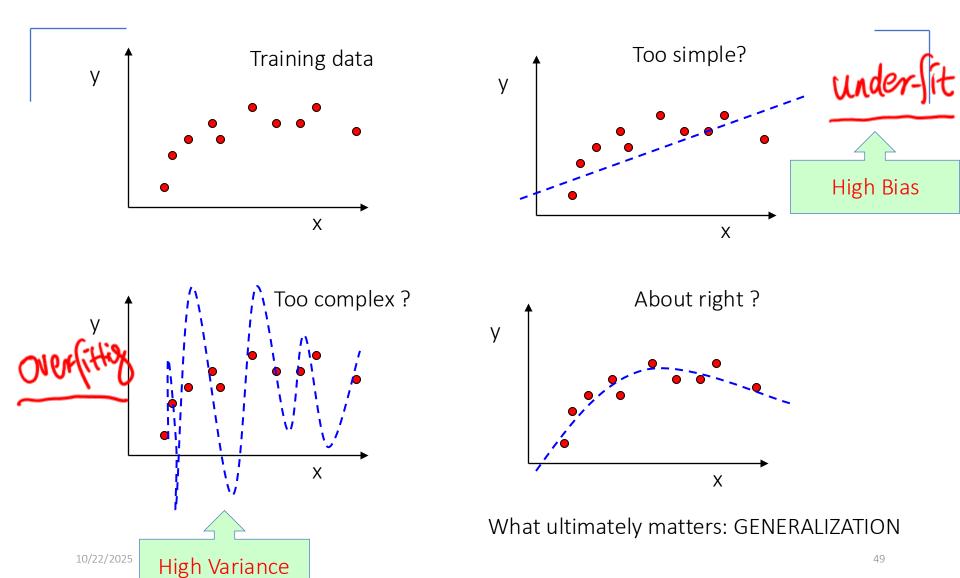
where ε is an error term of unmodeled effects or random noise

• Now assume that ε follows a Gaussian N(0, σ), then we have:

$$p(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

$$\text{RV} \quad \forall \mathbf{y} | \mathbf{x}_i; \theta \sim \mathbb{N} \left(\theta^T \mathbf{x}_i, \sigma\right)$$

L9: Complexity / Goodness of Fit / Generalization



Statistical Decision Theory (Extra)

- Random input vector: X
- Random output variable: Y
- Joint distribution: Pr(X,Y)
- Loss function L(Y, f(X))

Expected prediction error (EPE):

$$EPE(f) = E(L(Y, f(X))) = \Box L(y, f(x)) Pr(dx, dy)$$

$$e.g. = \Box (y - f(x))^{2} Pr(dx, dy)$$

e.g. Squared error loss (also called L2 loss)

One way to define generalization: by considering the joint population distribution

Decomposition of EPE

- When additive error model: $Y = f(X) + \epsilon, \ \epsilon \sim (0, \sigma^2)$
- Notations
 - ullet Output random variable: Y
 - True function: $f \rightarrow true$
 - Prediction estimator: $\hat{f} \rightarrow \mathcal{D} \rightarrow \hat{T}$

$$EPE(x) = E[(Y - \hat{f})^2 | X = x]$$

$$= E[((Y - f) + (f - \hat{f}))^2 | X = x]$$

$$= E[(Y - f)^2 | X = x] + E[(f - \hat{f})^2 | X = x]$$

$$= \sigma^2 + Var(\hat{f}) + Bias^2(\hat{f})$$

Irreducible / Bayes error

Bias-Variance Trade-off for EPE:

EPE $(x) = noise^2 + bias^2 + variance$

Unavoidable error

Error due to incorrect assumptions

Error due to variance of training samples

BIAS AND VARIANCE TRADE-OFF for Parameter Estimation

- θ : true value (normally unknown)
- $\widehat{\theta}$: estimator
- $\bar{\theta}$: = $E[\hat{\theta}]$ (mean, i.e. expectation of the estimator)
- Bias $E[(\bar{\theta} \theta)^2]$
 - measures accuracy or quality of the estimator
 - low bias implies on average we will accurately estimate true parameter from training data
- Variance $E[(\hat{\theta} \bar{\theta})^2]$
 - Measures precision or specificity of the estimator
 - Low variance implies the estimator does not change much as the training set varies

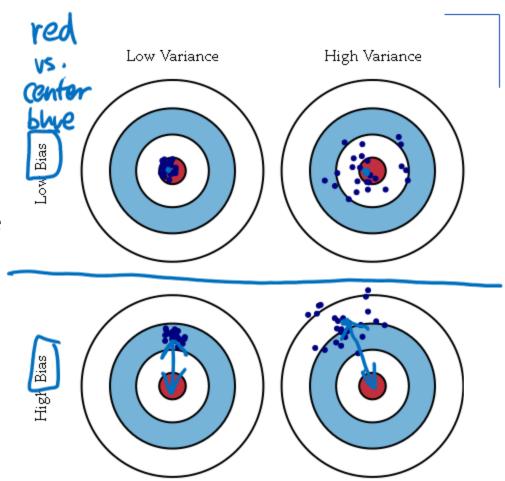
Model "bias" & Model "variance"

- Middle RED:
 - TRUE function
- Error due to bias:
 - How far off in general from the middle red

$$E[(\bar{\theta}-\theta)^2]$$

- Error due to variance:
 - How wildly the blue points spread

$$E[(\hat{\theta} - \bar{\theta})^2]$$



need to make assumptions that are able to generalize

- Underfitting: model is too "simple" to represent all the relevant characteristics
 - High bias and low variance
 - High training error and high test error
- Overfitting: model is too "complex" and fits irrelevant characteristics (noise) in the data
 - Low bias and high variance
 - Low training error and high test error

Bias Variance Tradeoff

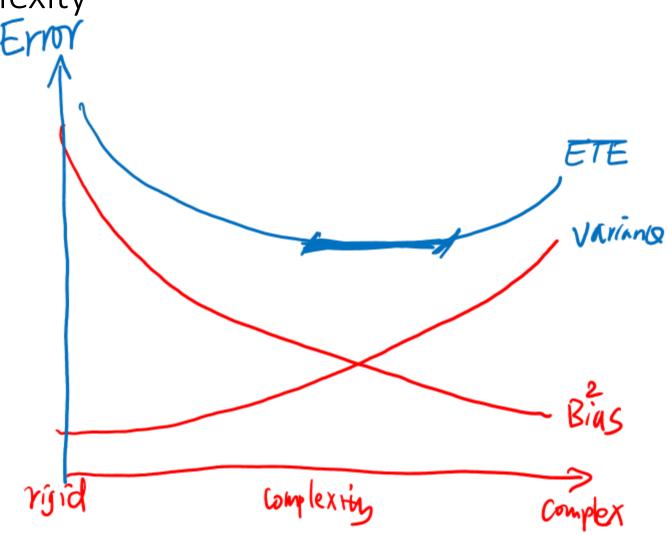
- •(1) Randomness of Training Sets
- (2) Training error can always be reduced when increasing model complexity
- •(3) Randomness in the Testing Error!!!
- (4) Cross Validation Error as good approximation for Expected Test error -- good appx of generalization

Review:

One important Control of Bias Variance Tradeoff

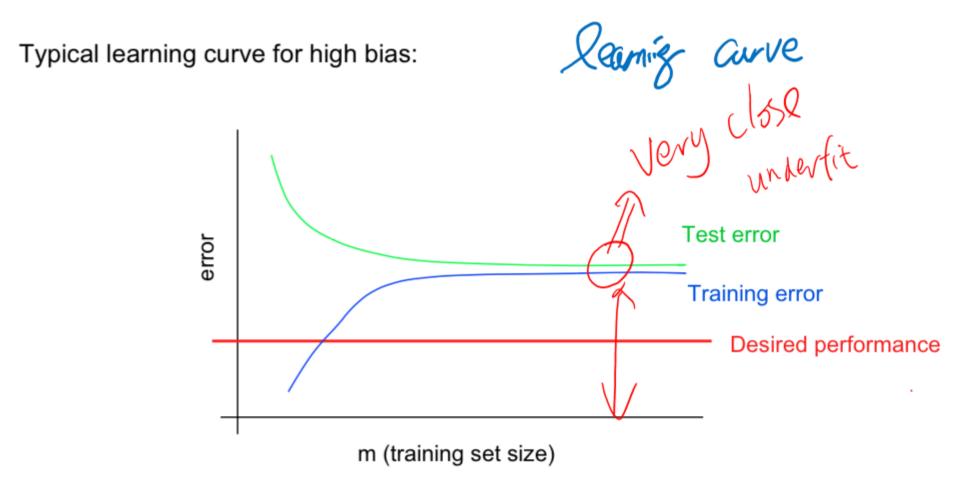
→ Model Complexity

- bias decrease with model gets more complex;
- Variance increase with bigger model capacity
- Sum of Bias^2+Variance



Another important Control of Bias Variance Tradeoff

→ Training Size (Extra)

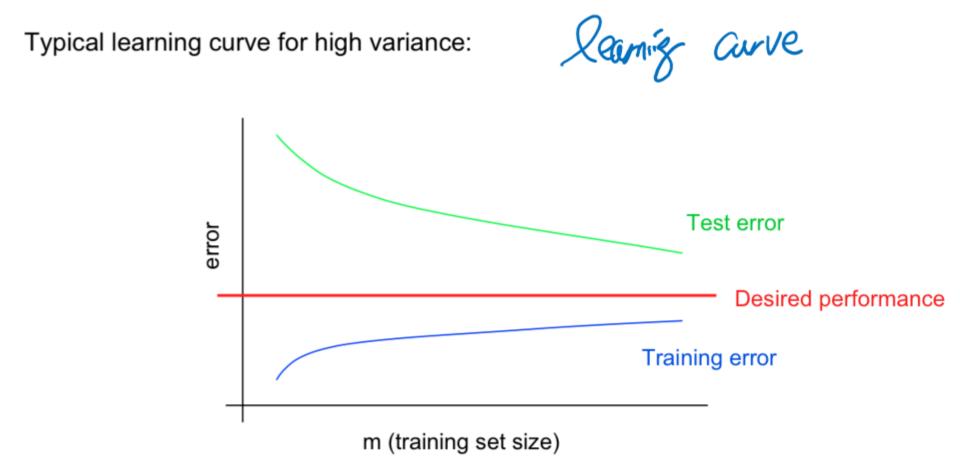


- Even training error is unacceptably high.
- Small gap between training and test error.

High training error and high test error

Another important Control of Bias Variance Tradeoff

→ Training Size (Extra)



How to reduce Model High Variance?

- Choose a simpler classifier
- Regularize the parameters
- Get more training data

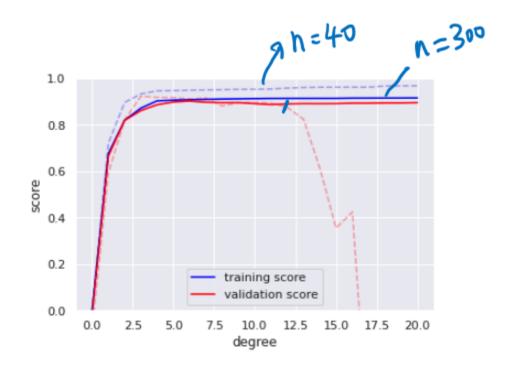


- Try feature engineering
- Try multiple models and then use all as ensemble

Take Away: Three types of plots

- (1) Sanity check (S)GD type Optimization
 - Train / Vali Loss vs. Epochs to help you
 - https://scikit-learn.org/stable/auto_examples/linear_model/plot_sgd_early_stopping.html#sphx-glr-auto-examples-linear-model-plot-sgd-early-stopping-py
- (2) Sanity check hyperparameter tuning (validation curve)
 - Train / Vali Loss vs. hyperparameter Values
 from sklearn.model_selection import validation curve
- (3) Sanity check if your current model overfits or underfits
 - Train / Vali Loss vs. Varying Size of Training (learning curve)
 - https://scikit-learn.org/stable/auto_examples/model_selection/plot_learning_curve.html#sphx-glr-auto-examples-model-selection-plot-learning-curve-py

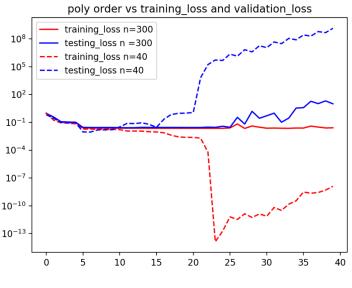
I will Code run: https://github.com/qiyanjun/2025Fall-UVA-CS- MachineLearningDeep/blob/main/notebook/L9-LearningCurves.ipynb



(1) Validation curve

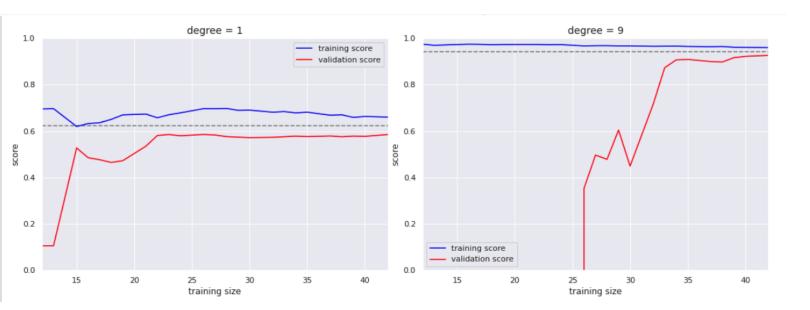
By scikitlearn Validation_curve function (normalize all metrics to positive range

https://scikit-learn.org/stable/modules/model_evaluation.html#the-scoring-parameter-defining-model-evaluation-rules

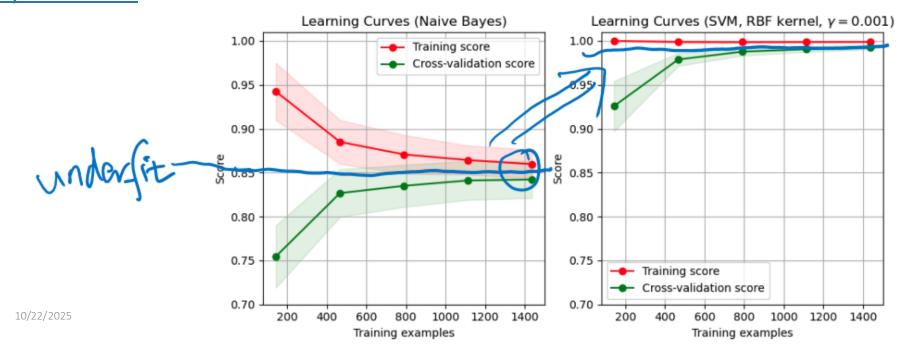


(1) Validation curve

By our HW2 (more close to modern deep learning library style)



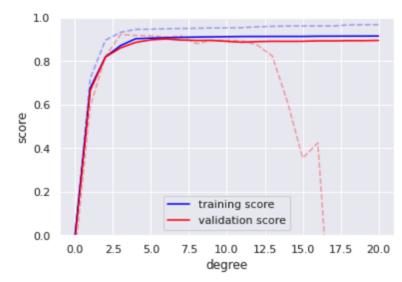
(1) Learning Curves for polynomial regression (up) and classification (down) / by scikitlearn



```
X2, y2 = make data(200)
degree = np.arange(200)
train score2, val score2 = validation curve(PolynomialR
                                                'polynomial
plt.plot(degree, np.median(train score2, 1), color='blu
plt.plot(degree, np.median(val score2, 1), color='red',
plt.legend(loc='lower center')
plt.ylim(0, 1)
plt.xlabel('degree')
plt.ylabel('score');
   1.0
   0.8
                                                           C→
   0.6
 score
   0.2
                        training score
                         validation score
```

0.0

25



Interesting Relation between

125

degree

150

the right range of model complexity

175

the number of training points

Is the bias-variance trade off dependent on the number of samples? (EXTRA)



In the usual application of linear regression, your coefficient estimators are unbiased so sample size is irrelevant. But more generally, you can have bias that is a function of sample size as in the case of the variance estimator obtained from applying the population variance formula to a sample (sum of squares divided by n).....

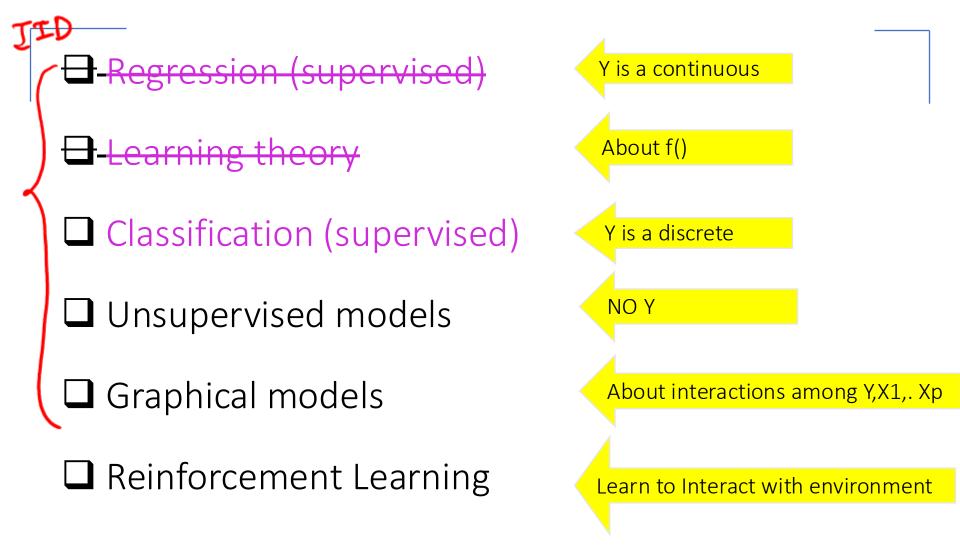
... the bias and variance for an estimator are generally a decreasing function of training size n. Dealing with this is a core topic in nonparametric statistics. For nonparametric methods with tuning parameters a very standard practice is to theoretically derive rates of convergence (as sample size goes to infinity) of the bias and variance as a function of the tuning parameter, and then you find the optimal (in terms of MSE) rate of convergence of the tuning parameter by balancing the rates of the bias and variance. Then you get asymptotic results of your estimator with the tuning parameter converging at that particular rate. Ideally you also provide a databased method of choosing the tuning parameter (since simply setting the tuning parameter to some fixed function of sample size could have poor finite sample performance), and then show that the tuning parameter chosen this way attains the optimal rate.

10/16

Agenda

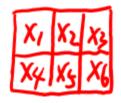
- Going over HW2 Solution
- A tutorial talk on Huggingface.co

Course Content Plan → Regarding Tasks



Course Content Plan
Regarding Data

- ☐ Tabular / Matrix
- atrix Si
- ☐ 2D Grid Structured: Imaging

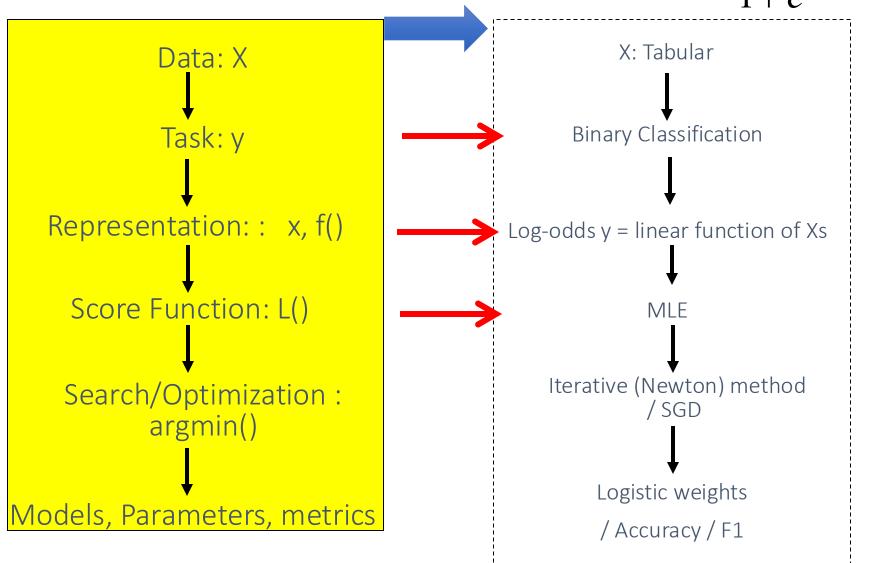


- ☐ 1D Sequential Structured: Text
- ☐ Graph Structured (Relational)
- ☐ Set Structured / 3D /

10/22/2025 69

Today: Logistic Regression Classifier

$$P(y=1|x) = \frac{e^{\beta_0 + \beta^T x}}{1 + e^{\beta_0 + \beta^T x}}$$



Bayes Classifiers – Predict via MAP Rule

Task: Classify a new instance X: $X = \langle X_1, X_2, ..., X_p \rangle$

based on:

$$c_{MAP} = \underset{c_j \square C}{\operatorname{argmax}} P(c_j \mid x_1, x_2, \dots, x_p)$$
MAP Rule

MAP = Maximum Aposteriori Probability

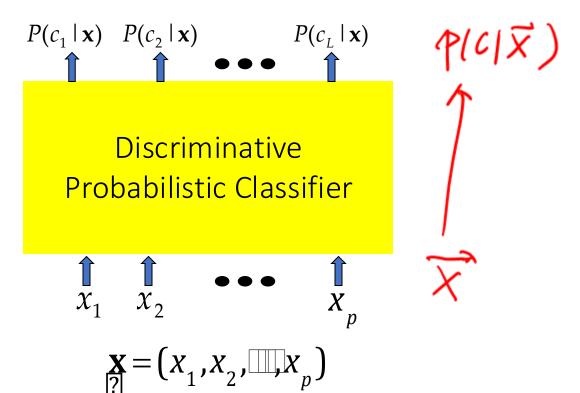
Our Whole Section 2:

$$X \longrightarrow C$$

φ(c|x)

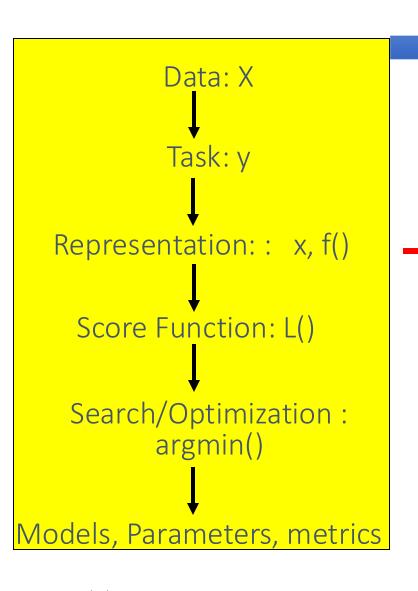
Discriminative Classifiers

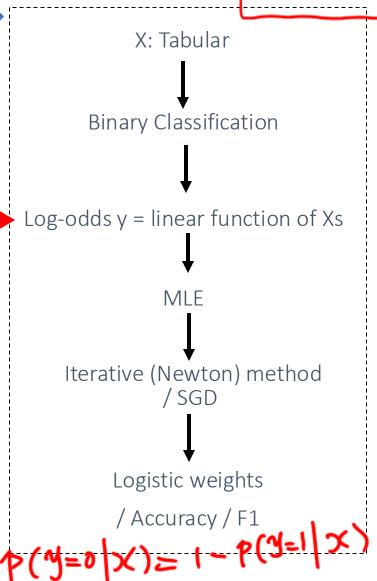
$$\underset{c \square C}{\operatorname{arg\,max}} P(c \mid \mathbf{X}), \mathbf{MC} = \{c_1, \mathbf{LC}\}$$



Today: Logistic Regression Classifier

$$P(y=1|x) = \frac{e^{\beta_0 + \beta^T x}}{1 + e^{\beta_0 + \beta^T x}}$$





Logistic Regression p(y|x)

$$\ln \left[\frac{P(y|x)}{P(T_{ai}(x))} \right] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$$|P(y|x)| = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$$|P(y|x)| = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$$|P(y|x)| = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$$= P(y|x) = \frac{e^{\rho_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}} = \frac{1}{1 + e^{-(\beta_0 + \beta^T X)}}$$

10/22/2025

View IV: Logistic Regression models a linear classification boundary!

$$\frac{p(y=0|x) = p(y=1|x)}{\log\left(\frac{p(y=1|x)}{p(y=0|x)}\right) = \beta \lambda} = \log(1) = 0$$
Decision Boundary
$$\frac{p(y=0|x) = p(y=1|x)}{p(y=0|x)} = 1$$

SVM

10/22/2025

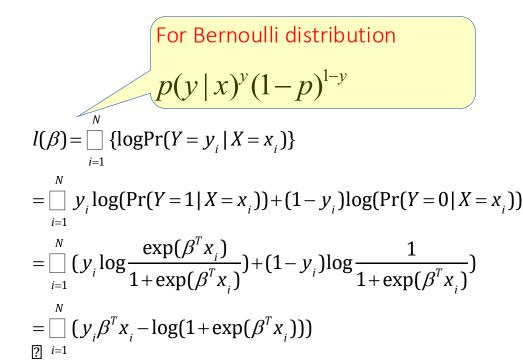
Summary: MLE for Logistic Regression Training

Let's fit the logistic regression model for K=2, i.e., number of classes is 2

Training set: (x_i, y_i) , i=1,...,N

(conditional) Log-likelihood:



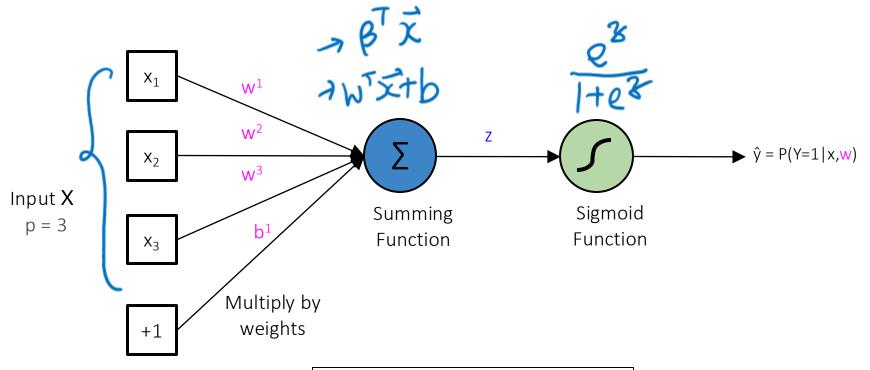


 x_i are (p+1)-dimensional input vector with leading entry 1 \beta is a (p+1)-dimensional vector

We want to maximize the log-likelihood in order to estimate \beta

See Extra Slides How to used Newton-Raphson optimization

One "Neuron": Block View of Logistic Regression



$$z = w^{T} \cdot x + b$$

$$y = sigmoid(z) = \frac{e^{z}}{1 + e^{z}}$$

10/21

Agenda

- •HW3 is due
- Please select your Project's Shark Tank
 Sessions ASAP

- •Today:
 - Review MLP / DNN / CNN / PCA / Word Embedding, Transformer
 - Quiz 8 Today

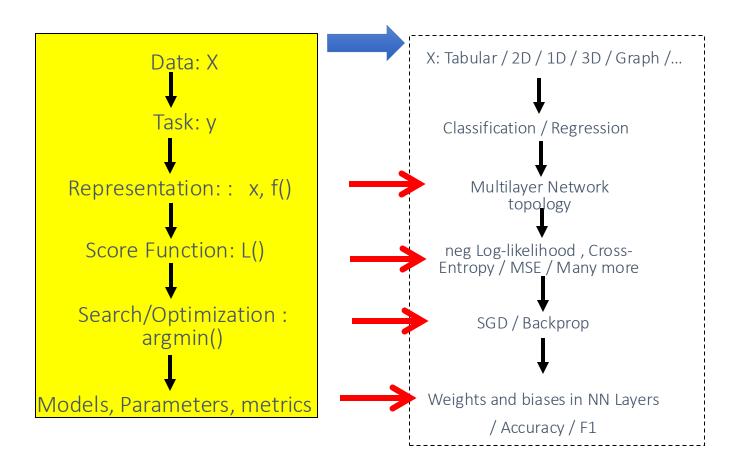
Takeaway: Logistic Regression Classifier

- View I: logit(y) as linear of Xs
- View II: model Y as Bernoulli with p(y=1|x) as p(Head)
- View III: S" shape function compress to [0,1]
- View IV: models a linear classification boundary!
- View V: Two stages: summation + sigmoid

10/22/2025

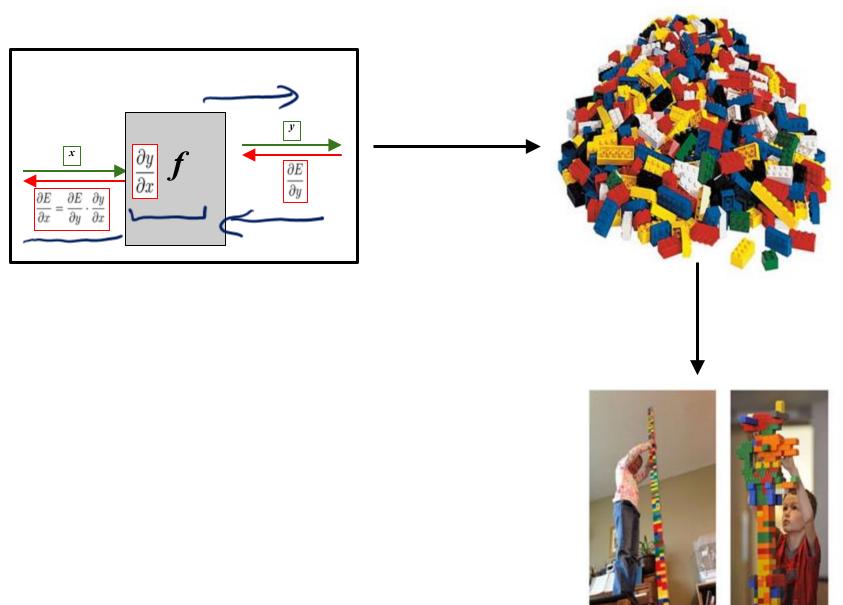
Lecture 12: Neural Network (NN) and More: BackProp

Today: Basic Neural Network Models

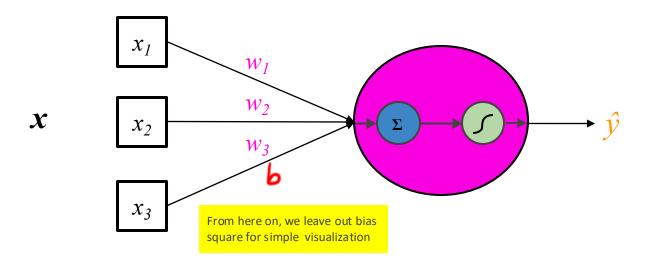


10/22/2025

Building Deep Neural Nets

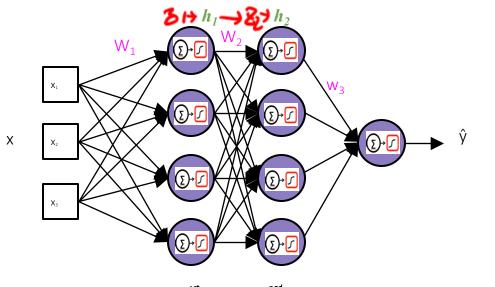


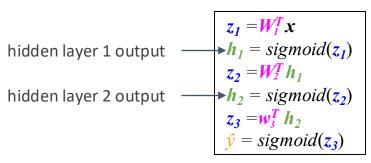
Neuron Representation



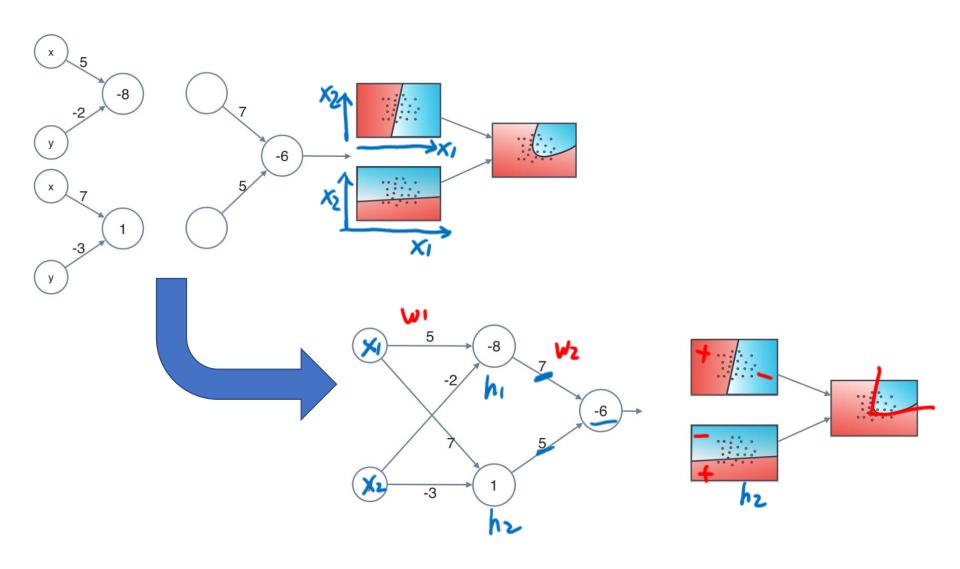
The linear transformation and nonlinearity together is typically considered a single neuron

Multi-Layer Perceptron (MLP)- (Feed-Forward NN)





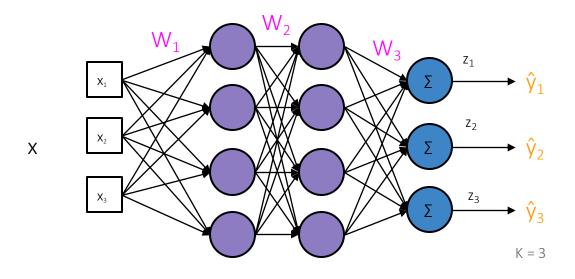
$$\times \overrightarrow{w}$$
 $\delta_1 \rightarrow h_1 \xrightarrow{W_2} \delta_2 \rightarrow h_2 \xrightarrow{\widetilde{W}} \widetilde{Y}$

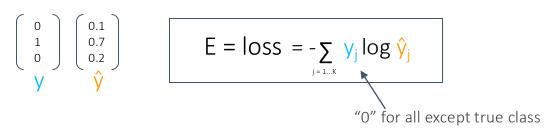


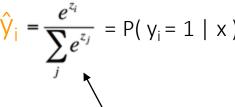
10/22/2025

Recap: Multi-Class Classification Loss

Cross Entropy Loss



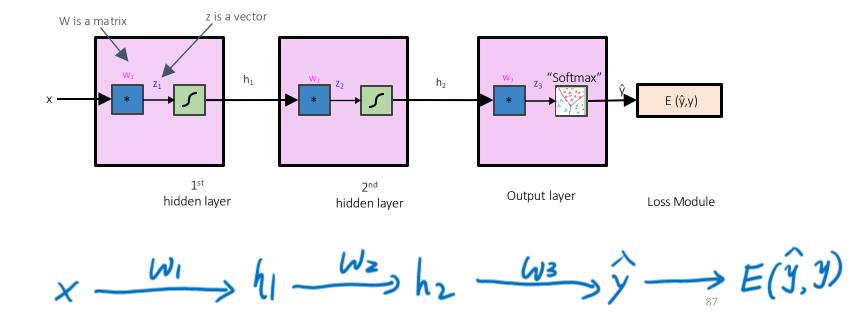




"Softmax" function.

Normalizing function which converts each class output to a probability.

e.g., "Block View" of multi-layered multi-class NN



Extra

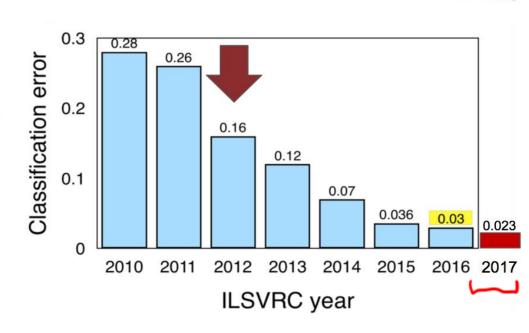
argmin_w $\{f_4(f_3(f_2(f_1())))\}$			Local Gradients= ∂ Output / ∂ Input				
	Input	Output	Local Gradients-O Output /Omput				
	ار ، ۲۲ سار . الار الار الار الار الار الار الار	31,82	$\frac{\partial \mathcal{E}_1}{\partial \mathcal{N}_1} = \mathcal{N}_1$				
f_{2} $h_{1} = \frac{exp(z_{1})}{T + exp(z_{1})}$ $h_{2} = \frac{exp(z_{2})}{T + exp(z_{2})}$	31,32	hishz	3h1 = h1 (1-h1)				
$\hat{\mathbf{y}} = \mathbf{h}_1 \mathbf{w}_5 + \mathbf{h}_2 \mathbf{w}_6 + \mathbf{b}_3$	ws, hishz	প্	an/ahi= Ws				
$f_4 E = (y - \hat{y})^2$	Ŷ	losse	$\partial E/\partial \hat{y} = -2(y-\hat{y})$				
$\frac{\partial E}{\partial M} = \frac{\partial F_{u}}{\partial N_{1}} = \frac{\partial f_{u}}{\partial f_{3}} \frac{\partial f_{2}}{\partial f_{2}} \frac{\partial f_{1}}{\partial M_{1}}$ $= -2(9-9) \frac{\partial (h_{1}W_{5} + h_{2}W_{6} + b_{3})}{\partial W_{1}}$ $= -2(9-9) (W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial N_{1}})$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} = -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} \frac{\partial h_{1}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} = -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} \frac{\partial h_{1}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} = -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} \frac{\partial h_{1}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{2}}{\partial W_{1}} + W_{6} \frac{\partial h_{1}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \frac{\partial h_{1}}{\partial W_{1}}$ $= -2(9-9) W_{5} \frac{\partial h_{1}}{\partial W_{1}} + W_{6} \partial h_$							

Lecture 13: Supervised Image Classification and Convolutional Neural Networks

ImageNet Challenge

Arch

- 2010-11: hand-crafted computer vision pipelines
- 2012-2016: ConvNets
 - 2012: AlexNet
 - major deep learning success
 - 2013: ZFNet
 - improvements over AlexNet
 - 0 2014
 - VGGNet: deeper, simpler
 - InceptionNet: deeper, faster
 - 0 2015
 - ResNet: even deeper
 - 2016
 - ensembled networks
 - 0 2017
 - Squeeze and Excitation Network



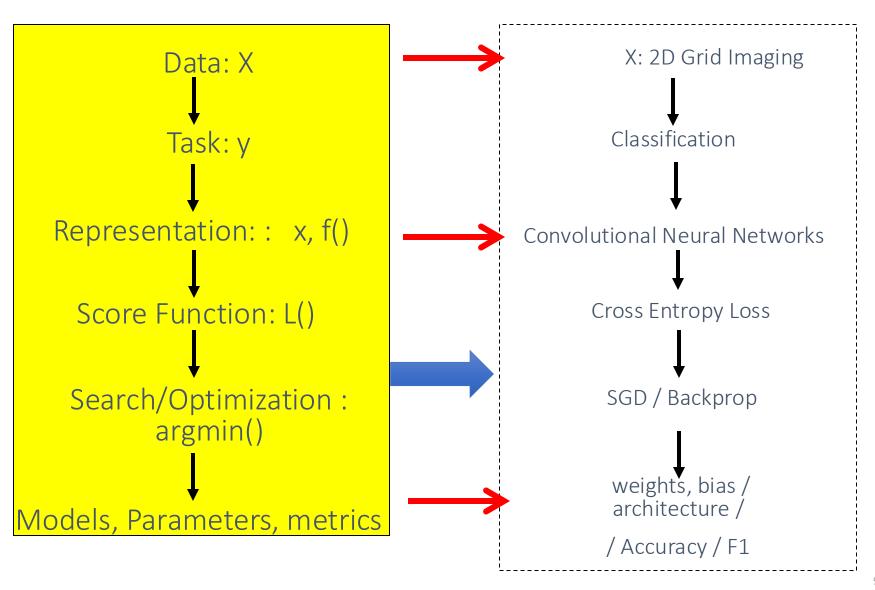
Lecture 13: Supervised Image Classification and Convolutional Neural Networks

		,		
Metric	Formula	Interpretation		
Accuracy	$\frac{\mathrm{TP} + \mathrm{TN}}{\mathrm{TP} + \mathrm{TN} + \mathrm{FP} + \mathrm{FN}}$	Overall performance of model		
Precision	$\frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FP}}$	How accurate the positive predictions are		
Recall Sensitivity	$\frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FN}}$	Coverage of actual positive sample		
Specificity	$\frac{\mathrm{TN}}{\mathrm{TN} + \mathrm{FP}}$	Coverage of actual negative sample		
F1 score	$\frac{2\mathrm{TP}}{2\mathrm{TP} + \mathrm{FP} + \mathrm{FN}}$	actual ed classes		
		predicted+ TP FP		

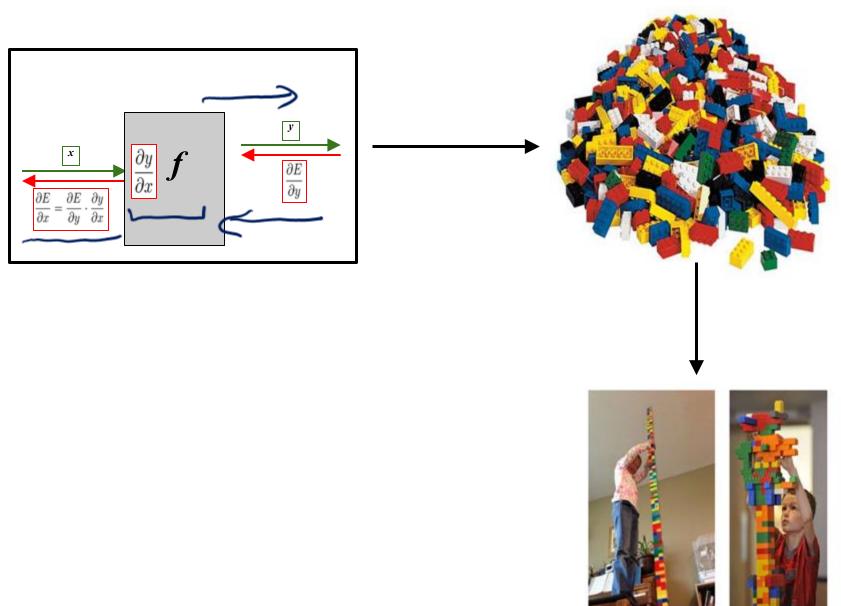
10/22/2025

Lecture 13: Supervised Image Classification and Convolutional Neural Networks

Today: Convolutional Network Models on 2D Grid / Image



Building Deep Neural Nets



CNN models Locality and Translation Invariance

Important Block: Convolutional Neural Networks (CNN)

- Prof. Yann LeCun invented CNN in 1998
- First NN successfully trained with many layers







The bird occupies a local area and looks the same in different parts of an image.

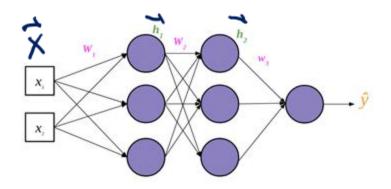
We should construct neural nets which exploit these properties!

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, Gradient-based learning applied to document recognition, Proceedings of the IEEE 86(11): 2278–2324, 1998.

TF Keras Sample Code

https://www.kaggle.com/code/shawon10/covid-19-diagnosis-from-images-using-densenet121

Pytorch Sample Code



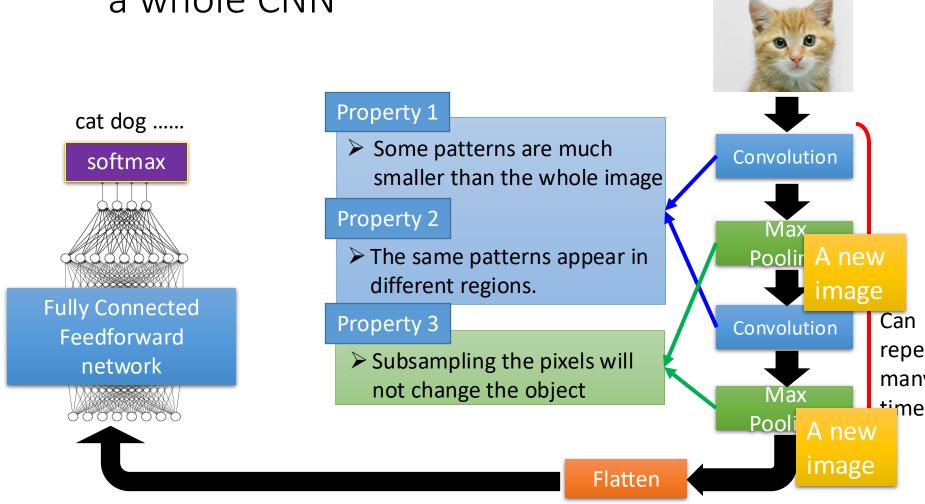
```
import torch.nn as nn
import torch.nn.functional as F

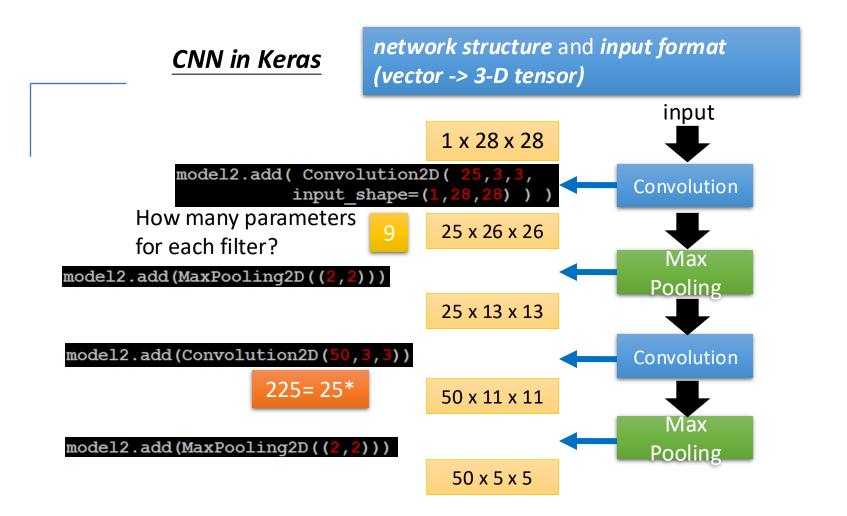
class ThreeLayerNet(torch.nn.Module):
    def __init__(self, d_in, d_hidden, d_out):
        super().__init__()
        self.W1 = nn.Linear(d_in,d_hidden)
        self.W2 = nn.Linear(d_hidden,d_hidden)
        self.w3 = nn.Linear(d_hidden,d_out)
        self.nonlinear = nn.Sigmoid()

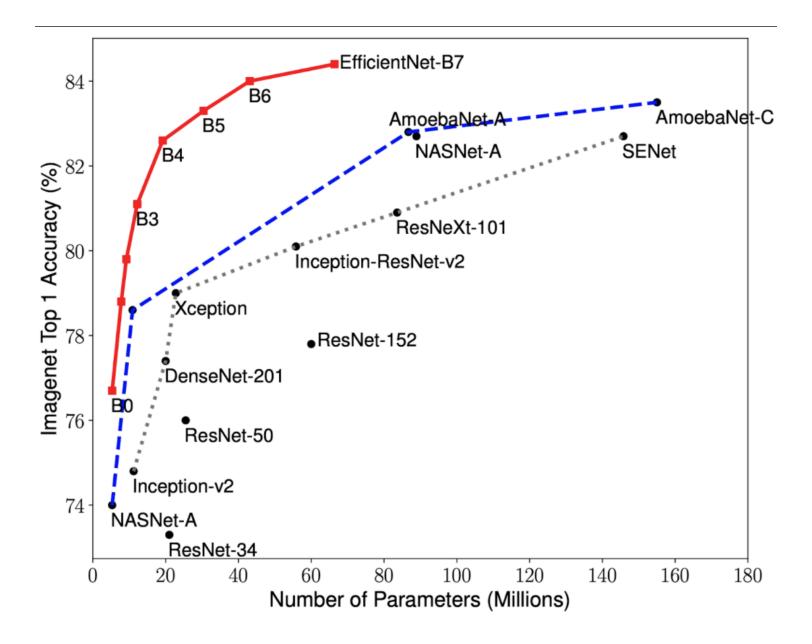
def forward(self, x):
    h1 = self.nonlinear(self.W1(x))
    h2 = self.nonlinear(self.W2(h1))
    y_hat = self.nonlinear(self.w3(h2))
    return y_hat

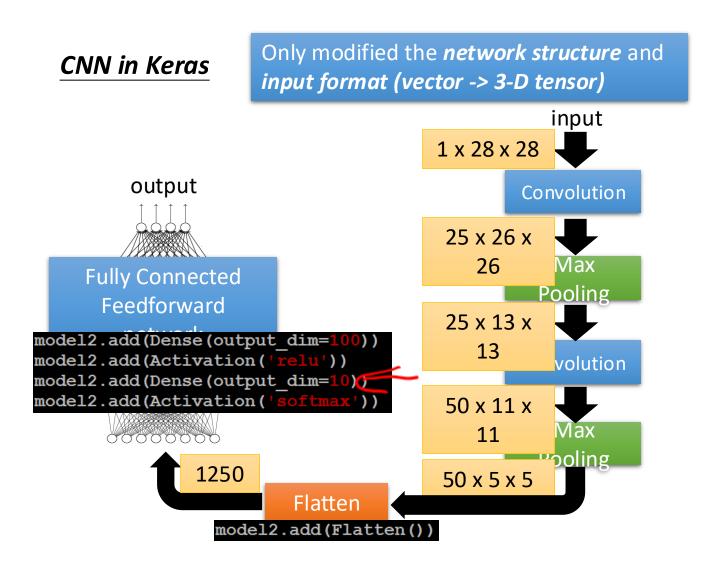
model = ThreeLayerNet(2,3,1)
```

a whole CNN









Lecture 14: Dimension Reduction

Today: Dimensionality Reduction (Two Ways)

Feature selection: chooses a subset of the original features.

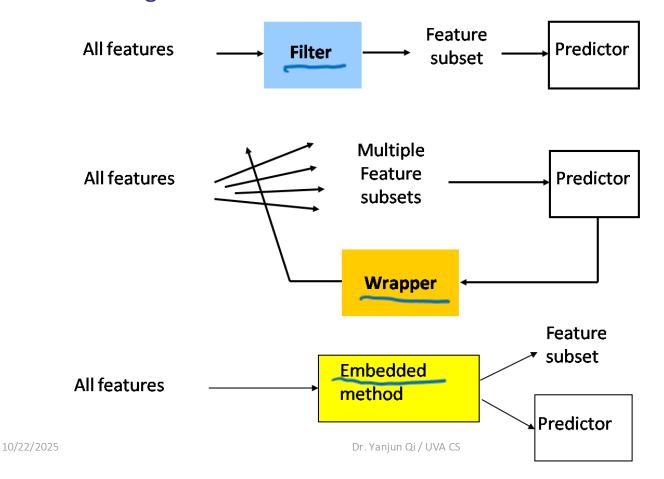


$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \longrightarrow \mathbf{x}' = \begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{kK} \end{bmatrix}$$

K<<N

Summary: Feature Selection => filters vs. wrappers vs. embedding

Main goal: rank subsets of useful features



(I) Filtering: (many choices)

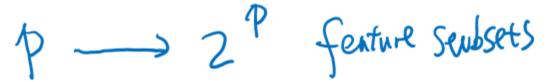
Method	X Y	mments		
Name $ Formula B M C B M C $				
Bayesian accuracy Balanced accuracy		_	d, rescaled Bayesian relevance Eq. 3.2. icity; used for unbalanced dataset,	
Bi-normal separation F-measure Odds ratio	Eq. 3.5 + s + s Eq. 3.7 + s + s Eq. 3.6 + s + s	ed in information retrieval.	, popular in information retrieval.	
Means separation T-statistics Pearson correlation Group correlation χ^2 Relief Separability Split Value	Eq. 3.13 + i + i i + i Eq. 3.8 + s + s	arson's coefficient for subset o sults depend on the number o	tion. est Eq. 3.12 , or a permutation test. f features. f samples m . is for a simplified version ReliefX,	
Kolmogorov distance Bayesian measure Kullback-Leibler divergence Jeffreys-Matusita distance Value Difference Metric	Eq. 3.16 + s + + s Eq. 3.20 + s + + s	ference between joint and prome as Vajda entropy Eq. 3.23 uivalent to mutual informationally used but worth trying. ed for symbolic data in similally symbolic feature-feature cort	and Gini Eq. 3.39. n. rity-based methods,	
Mutual Information V Information Gain Ratio V Symmetrical Uncertainty J-measure Weight of evidence MDL 10/22/2025	Eq. 3.32 + s + + s Eq. 3.35 + s + + s	w bias for multivalued feature easures information provided b	ure entropy, stable evaluation. s. by a logical rule. Guyon-Flisseeff IMLR 2004	

(2) Wrapper: Feature Subset Selection

Wrapper Methods

- Learner is considered a black-box
- Interface of the black-box is used to score subsets of variables according to the predictive power of the learner when using the subsets.
- Results vary for different learners

(b). Search: even more search strategies for selecting feature subset



- Forward selection or backward elimination.
- Beam search: keep k best path at each step.
- **GSFS:** generalized sequential forward selection when (n-k) features are left try all subsets of g features. More trainings at each step, but fewer steps.
- PTA(I,r): plus I, take away r at each step, run SFS I times then SBS r times.
- Floating search: One step of SFS (resp. SBS), then SBS (resp. SFS) as long as we find better subsets than those of the same size obtained so far.

(3) Embedded

•Embedding approach:

uses a predictor to build a (single) model with a subset of features that are internally selected.

lasso elastiNet

Today: Dimensionality Reduction (Two Ways)

Feature extraction: finds a set of new features (i.e., through some mapping f()) from the existing features.

Feature selection: chooses a subset of the original features.



The mapping f() could be linear or non-linear

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \xrightarrow{\mathbf{f()}} \mathbf{h} = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_K \end{bmatrix}$$

$$\mathbf{K} \leq \mathbf{N}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \longrightarrow \mathbf{x}' = \begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{kK} \end{bmatrix}$$

K<<N

Feature Extraction (linear or nonlinear)

- Linear combinations are particularly attractive because they are simpler to compute and analytically tractable.
- Given $x \in \mathbb{R}^p$, find an N x K matrix U such that:

$$y = U^Tx \in R^K$$
 where K

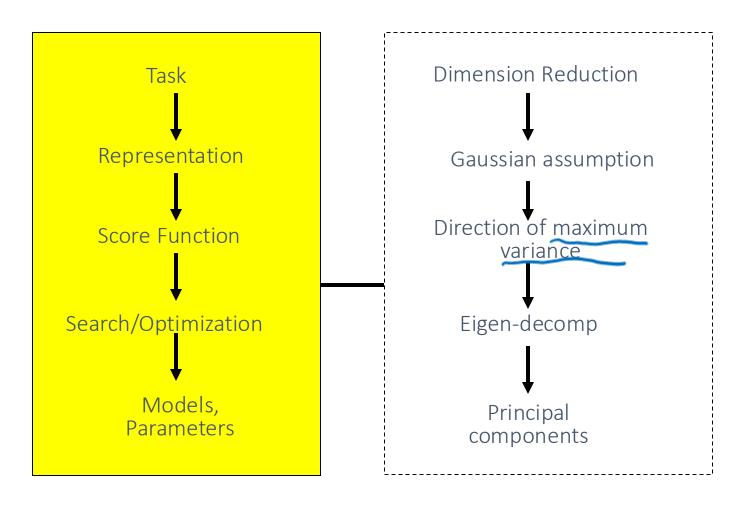
This is a projection from the N-dimensional space to a K-dimensional space.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \longrightarrow \mathbf{h} = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_K \end{bmatrix}$$

Feature Extraction (cont'd)

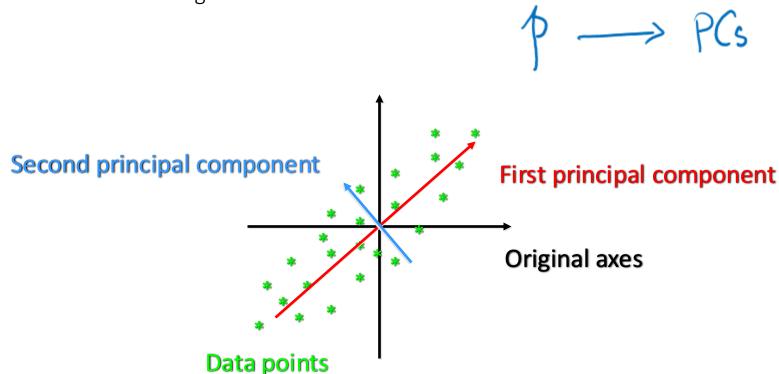
- Commonly used linear feature extraction methods:
 - Principal Components Analysis (PCA): Seeks a projection that preserves as much information in the data as possible.
 - Linear Discriminant Analysis (LDA): Seeks a projection that **best** discriminates the data.
- Recent nonlinear feature extraction methods:
 - Like Word Embedding / Autoencoder / ...

Principal Component Analysis

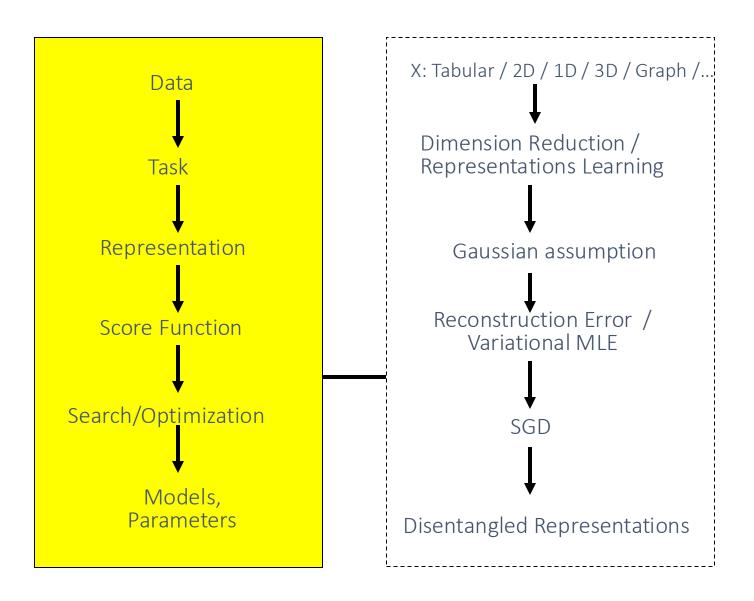


How does PCA work? Explaining Variance

- Each PC always explains some proportion of the total variance in the data. Between them they explain everything
 - PC1 always explains the most
 - PC2 is the next highest etc. etc.

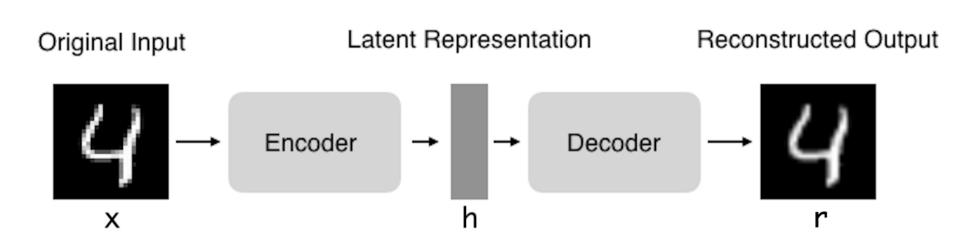


Auto Encoder

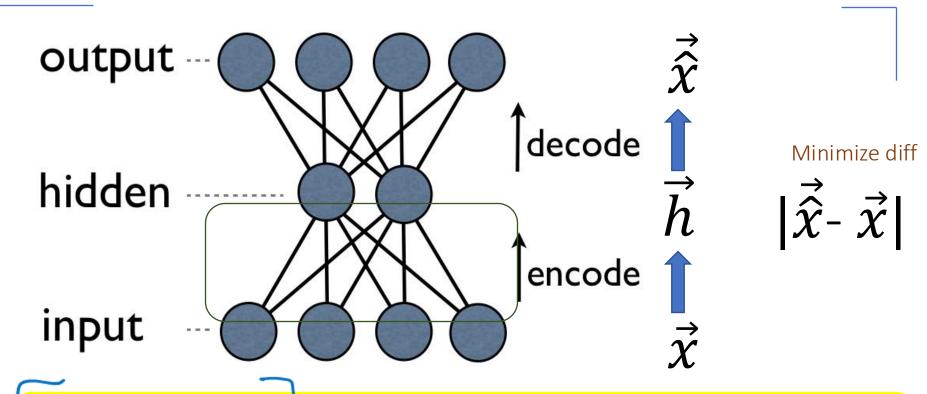


Autoencoders: structure

- Encoder: compress input into a latent-space of usually smaller dimension. h = f(x)
- Decoder: reconstruct input from the latent space. r = g(f(x)) with r as close to x as possible



an auto-encoder-decoder is trained to reproduce the input



Reconstruction Loss: force the 'hidden layer' units to become good / reliable feature detectors

10/30/19 Yanjun Qi / UVA CS 112

10/28

S3: Lecture 18:

Deep Neural Networks for Natural Language Processing

How to Represent A Word in DNN: Feature Extraction / Embedding

- Basic approach "one hot vector"
 - Binary vector
 - Length = | vocab |
 - 1 in the position of the word id, the rest are 0
 - However, does not represent word meaning
 - Extremely high dimensional (there are over 200K words in the English language)
 - Extremely sparse

• Solution:

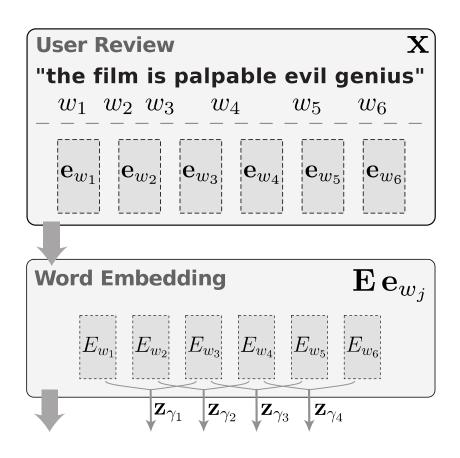
Distributional Word Embedding Vectors

Popular word embeddings

- GloVe (Global Vectors)
 - o Pennington et al., 2014
- fasttext
 - Bojanowski et al., 2017

However, Natural language is

- Variable-length
- Composition of multiple words
- Word meaning is contextual
 - Elmo
 - Peters, 2018
 - BERT
 - o Devlin et al., 2018



S3: Lecture **18:**

Deep Neural Networks for Natural Language Processing