

Lecture 2: Regression Analysis

(closed book, simple calculator)

Confidence Interval

Provide us w/ range of values that we believe, with % level of confidence, contain a population parameter

$$\Pr\left(\bar{X} - z_{(1-\frac{\alpha}{2})} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{(1-\frac{\alpha}{2})} \frac{\sigma}{\sqrt{n}}\right) \quad \alpha = \% \text{ level confidence}$$

Larger sample = smaller interval

Lower α = smaller interval, larger variance = larger interval

use z if know population σ , use t if know sample s

$$Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \quad t = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \quad \text{df} = n-1 \rightarrow \text{use in the t-dist. graph}$$

Hypothesis Testing

1. State null hypothesis H_0

and Alternative hypo. H_1

2. Choose α - level

or Calculate p-value

if $p \leq \alpha$ reject H_0

$p > \alpha$ cannot reject H_0 (not accept H_0 tho)

H_0 H_1 Reject H_0 if

$\mu = \mu_0$ $\mu \neq \mu_1$ $|z| \geq z_{1-\alpha/2}$ or $|t| \geq t_{1-\alpha/2, n-1}$ (2 sided test)

$\mu \leq \mu_0$ $\mu > \mu_0$ $z \geq z_{1-\alpha}$ or $t \geq t_{1-\alpha, n-1}$ (one sided test)

$\mu \geq \mu_0$ $\mu < \mu_0$ $z \leq z_{\alpha} = -z_{1-\alpha}$ or $t \leq t_{\alpha, n-1}$ (one sided test)

Type I Error : False Positive - Reject H_0 when H_0 is actually true

Type II Error : False Negative - Fail to reject H_0 when H_1 is true

P-Value and t-Test explanation

The hypothesis that there's no difference between 2 measured is called Null Hypothesis (H_0). In T-Test, we check with 95% (more or less) confidence interval which is the middle curve part of the normal distribution.

If the **t test is higher** than the t value of H_0 at 95% confident. Then that mean that the new distribution is not inside the 95% of H_0 and there's enough difference between the 2 distribution, so we can **reject the H_0** (it does not exist within H_0 area).

If **it's lower**, then it's inside the 95% H_0 area and the difference is not significant enough to reject the H_0 (that there's no difference) so we **fail to reject H_0** .

P-Value work similarly. If **p-value is large**, then the difference between 2 distribution is "**not significant**" enough, meaning there's no significant diff between default and the one we're measuring. The measured data is close enough to the distribution. If we create the distribution on the measured data, it will be not be significantly difference from the first distribution so **there's high probability that we can observe the value under H_0 area**. (and **fail to reject H_0**)

If the **p-value is small**, then the difference between 2 distribution is "**significant**", meaning the measured data is far away from the distribution and other distribution would do a better job explaining the data. If we create distribution on measured data, the distribution will be significant away from the 1st distribution and **there's low probability that we will observe the value under H_0 area** (and **reject the H_0**)

Alternative Hypothesis is **another hypothesis that is true if null hypothesis is false**. It's mostly the opposite of null hypothesis (In this case, there's a difference between 2 measured)

t-Test Type

Single Sample - test whether \bar{x} is significantly diff from a pre-existing value

Paired sample - test relationship btw 2 linked sample eg. mean obtain in 2 condition by a single group

Independent sample - test relationship btw 2 indep. problem

Paired t-Test w/ 2 Paired Sample

$$H_0: \mu_d = \mu_1 - \mu_2 = D_0$$

Test statistic $t = \frac{\bar{d} - D_0}{\frac{s}{\sqrt{n}}}$ where \bar{d} = mean of the diff btw 2 sample test

$$\begin{aligned} H_1: \mu_d &\neq D_0 & \text{Reject } H_0 \text{ if: } |t| \geq t_{1-\frac{\alpha}{2}, n-1} \\ &: \mu_d > D_0 &: t \geq t_{1-\alpha, n-1} \\ &: \mu_d < D_0 &: t \leq -t_{\alpha, n-1} \end{aligned}$$

Independent Sample - 2 indep. sample

$H_0: m_1 = m_2$

Welch's t-Test

$H_1: m_1 \neq m_2$

2-sample unpaired t-test w/ equal or not equal sample size, assuming unequal variance

$$df = n_1 + n_2 - 2$$

$$t = \frac{\bar{x}_1 - \bar{x}_2}{S_{\bar{x}_1 - \bar{x}_2}}$$

$$\text{where } S_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Linear Regression

Relationship btw dependent variable (y = label/outcome) and indep. variable (x = all the input feature)

Simple Linear Regression model

$$y = \beta_0 + \beta_1 x + \varepsilon \quad \beta_0, \beta_1 \text{ are unknown coefficient}$$

β_0 = y-axis intercept β_1 = slope of the line ε = random error term

Ordinary Least Square (OLS) Estimators

Which line fit data the best? \rightarrow line that minimize Sum of Squared Residual (loss func)

$$\text{let } \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad RSS = \min \sum_i (y_i - \hat{y}_i)^2 = \sum_i (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x))^2$$

$$\text{then } \hat{\beta}_1 = \frac{\text{cov}(x, y)}{\text{var}(x)} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Note
 $\text{cov}(xy) = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{n-1}$
 $\text{var}(x) = \frac{\sum_i (x_i - \bar{x})^2}{n-1}$

Residual Sum of Square

$$RSS = \sum_i (y_i - \hat{y}_i)^2$$

$$ESS = \sum_i (\hat{y}_i - \bar{y})^2 \rightarrow \text{Explain sum of squared}$$

$$TSS = \sum_i (y_i - \bar{y})^2 = ESS + RSS = \text{Total sum of square}$$

R^2 measure proportion of variance in y that is explain by variance in x

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS} = \frac{ESS}{TSS} \quad (TSS = ESS + RSS)$$

$R^2 = 1 - \text{perfect match btw line and data}$

$R^2 = 0 - \text{no linear relationship btw } x \text{ & } y$

Multiple Linear Regression - $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_n x_n + \varepsilon$

$$y = \hat{x}\hat{\beta} + \varepsilon$$

OLS estimator - min error = min RSS = $\min \sum_i (y_i - \hat{y}_i)^2 = \min e^T e$

$$\frac{\partial RSS}{\partial \beta} = 0 \rightarrow \hat{\beta} = (X^T X)^{-1} X^T y \quad \text{and so} \quad \hat{y} = X \hat{\beta} = X (X^T X)^{-1} X^T y \quad (\text{normal equation})$$

Lecture 3 : Regression Diagnostics (BLUE)

Gauss-Markov Theorem

the OLS estimator is the best linear unbiased estimator (BLUE), iff the model has

1. Linearity - Linear relationship in parameters β

2. No multicollinearity of predictors $\rightarrow \text{Cov}(x_i x_j) = 0 \quad \forall i \neq j$

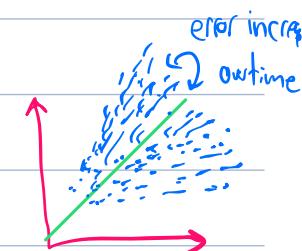
No linear dependency btw predictors

3. Homoscedasticity - The residual (error) exhibit constant variance
violate if variance increase overtime

4. No autocorrelation - no correlation btw i and j residual term (error)

5. Expected value of the residual vector (error) given X is 0 $E(\varepsilon | X) = 0$

exogeneity $\text{cov}(\varepsilon, X) = 0$



$$\text{unbiased} = E(\hat{\beta}_j) = \beta_j$$

best estimator = give lowest variance compare to other unbiased estimators (some biased can be better)

1. Linearity

For non linear \rightarrow reformulate model that make it linear again

- Polynomial regression - $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \varepsilon$ let $z = X^2 \rightarrow Y = \beta_0 + \beta_1 z + \beta_2 z^2 \rightarrow$ linear now

- Transform X and/or Y - $\log(Y) = \beta_0 + \beta_1 \log(X) + \varepsilon$

- Split data $X[X > X_k]$ (choose part of data that is linear)

2. No multicollinearity of predictors

Indep. var must not be linearly dependent. If 2 indep. var. are dependent \rightarrow omit one of them

- Check the rank of the matrix X must have full column rank $\text{rank}(X) = p$ $p = \# \text{ of col of } X$

High correlation btw indep. var leads to issue wrt. significance of predictors

Check for multicollinearity

Calculate the correlation coefficient btw every variable (both indep. and dep. var.)

Correlation = $\pm 1 \rightarrow$ high correlation Correlation = 0 no correlation

Variance Inflation Factor (VIF) = $\frac{1}{1-R_k^2}$ where R_k^2 = how much % of variance of predictor in question can be explained by other indep. var.

If VIF > 10, 90% of var has been explained by other. So we can remove that variable

If variable has low t-value or large p-value (non significant), then either

1. Variable is not related to response (y) \rightarrow small t , small VIF, small correlation w/ response

2. Relate to response (y) but strongly relate to other variable \rightarrow small t , large VIF, large correlation w/ response

(Look at example how variable change \rightarrow p.20-22)

Homoscedasticity

residual has a constant variance \rightarrow Homoscedasticity

if residual doesn't have constant variance \rightarrow Heteroscedasticity $\text{Var}(\varepsilon_i | x_{i1}, \dots, x_{ip})$ not constant

Test to check for homoscedasticity - Glejser Test, White Test

No Autocorrelation

No pattern should be detected when graphing residual over time

Modelling seasonality - use season instead of time b/c time has autocorrelation but season does not

>Create dummy variable which indicate season

\rightarrow have a time period to indicate which season is the data in

\rightarrow have 3 column for 3 season (Q_1, Q_2, Q_3), 1 if it's in one of the seasons, 0 if not

\rightarrow if all col is 0, then it's the last season \rightarrow if we model all 4 seasons \rightarrow multicollinearity

$$y = \beta_0 + \beta_1 T + \beta_2 Q_1 + \beta_3 Q_2 + \beta_4 Q_3$$

Expected value of the residual vector, given X , is 0 $E(\varepsilon | X) = 0$

Endogeneity $\rightarrow \text{corr}(\varepsilon_i, x_i) \neq 0 = E(\varepsilon_i | x_i) \neq 0$

Reason for endogeneity - omitted (confounding) variables

if True model: $y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon_i$

but our model use: $y_i = \beta_0 + \beta_1 x_1 + u_i$ then $u_i = \beta_2 x_2 + \varepsilon_i$

if x_1 and u_i are correlated and x_2 affects y , this lead to endogeneity

Technique to address endogeneity

Panel Data vs Cross Section Data

(not care about time) 

Cross-section data - observe many subjects at the same point of time, or who regard to diff in time

Panel Data Set - repeated observation on the same subject which make it possible to overcome omitted variable bias

↳ Unbalance - some subject haven't been record in some time period (data not full)

↳ Balance - every subject is record in everytime period

Model fixed Effects

Fixed effect - Treat γ_i as a constant for each subject (add to the intercept $\beta_0 + \gamma_i$)

- Idea: γ represent everything we don't know about the subject (unobserved heterogeneity).

Get a good estimate for other covariate (other variable that is correlate) w/o knowing the unobserved heterogeneity

Have various estimator for fixed effect model $y_i = (\beta_0 + \gamma_i) + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \varepsilon_i$

Lecture 4: Logistic & Poisson Regression

In the OLS regression, predicted prob. of the linear model can be >1 or <0

In Logistic Regression $= P[Y|X]$ - constrains estimated prob. to lie btw 0 and 1

$$P(Y|X) = p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

$P(Y=y|X)$ is \approx non linear func. β_0 = constant \rightarrow move curve left/right
 β_1 = slope = steepness of curve

Odds and Logit

$$\text{odd} = \frac{p(x)}{1-p(x)} = e^{\beta_0 + \beta_1 x}$$

odd range from 0 to ∞ (based on $e^{\beta_0 + \beta_1 x}$)

Logits \rightarrow take "ln" on both side $\ln\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x \rightarrow$ linear prob btw 0 and 1

logits is now between $-\infty$ to ∞ (based on $\beta_0 + \beta_1 x$) and β_0 = y-intercept, β_1 = slope

Estimating Coefficients of Logistic Regression

Maximum Likelihood estimate (MLE) \rightarrow Assume $X \sim \text{Ber}(p)$ so $p(y_i) = p^{y_i} (1-p)^{1-y_i}$

Likelihood func $L = \prod_{i=1}^n p^{y_i} (1-p)^{1-y_i}$

$$P(Y_i=1|X) = p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \sigma(\beta_0 + \beta_1 x), \text{ so } L = \prod_{i=1}^n \sigma(\beta_0 + \beta_1 x)^{y_i} (1 - \sigma(\beta_0 + \beta_1 x))^{1-y_i}$$

find β that
maximize likelihood

$$\text{Log Likelihood } LL = \ln(L) = \sum_i y_i \ln p_i + (1-y_i) \ln(1-p_i) \rightarrow \max_{\beta} \sum_i y_i \ln(\sigma(\beta_0 + \beta_1 x)) + (1-y_i) \ln(1 - \sigma(\beta_0 + \beta_1 x))$$

Gradient Ascent - $x^{(k+1)} = x^{(k)} + \alpha \nabla f(x^{(k)})$ gradient \rightarrow direction of steepest ascent

Use it to maximize LL or minimize NLL

Goodness of Fit

Null model - model with only the intercept (only estimate 1 param)

Fitted model - model w/ p parameters and intercept term, so we have p+1 param

Null Deviance = $-2\ln(L(\text{null}))$ - How much is explained by a null model

Residual Deviance = $-2\ln(L(\text{fitted}))$ - How well the fitted model explain the data \rightarrow small value = explain well

McFadden R²

Similar to R^2 , but R^2 is for linear regression, McFadden R^2 is for logistic regression.

$$R^2_{\text{McFadden}} = 1 - \frac{LL(\text{fitted})}{LL(\text{null})}$$

If the fitted model explain much more than a null model, then $R^2 \rightarrow 1$
 If fitted model doesn't explain much $\rightarrow R^2 \rightarrow 0$ ($R^2 > 0.2$ is okay)

logit model is an example of Generalised Linear model (GLM)

GLM - general class of linear model made up of 3 component

1. Distribution for modelling Y (Gaussian, Binomial, Poisson)

2. Linear Prediction $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

3. Link function $g(\cdot)$

{ Identity link (use in normal linear regression) - $g(\mu) = \mu = X\beta$

{ Logit Link - $g(\mu) = \log(\frac{\mu}{1-\mu})$

{ Log Link - $g(\mu) = \log(\mu)$

Count variable

Count variable \rightarrow non negative integer

Can be model w/ OLS regression but it can yield negative predicted value and count variable is often skewed

Use Poisson model or Negative binomial

In poisson, coefficient β_i , when interpret need to be put in $e^{\beta_i} \rightarrow$ Variable x_i decrease/increase rate by e^{β_i} factor

Lecture 5: Naive Bayes

Prob. can be multiply

Naive Bayes - All feature are equally important and independent of each other

Conditional Prob and Bayes Rule

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad P(A|B)P(B) = P(B|A)P(A) \quad (\text{product rule})$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (\text{Bayes Rule}) \quad P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

$$P(B) = P(B|A)P(A) + P(B|\neg A)P(\neg A) = P(B \cap A) + P(B \cap \neg A)$$

R = Result of test D = disease

$$\begin{aligned} P(R=P|D=P) &= 0.98 \rightarrow P(R=\neg P|D=P) = 0.02 \\ P(R=\neg P|D=P) &= 0.97 \rightarrow P(R=P|D=\neg P) = 0.03 \\ P(D=P) &= 0.008 \quad P(D=\neg P) = 0.992 \end{aligned}$$

$$\begin{aligned} P(D=P|R=P) &= \frac{P(R=P|D=P)P(D=P)}{P(R=P)} \rightarrow P(R=\neg P) - P(R=P|D=\neg P)P(D=\neg P) \\ &= \frac{(0.98)(0.008)}{(0.97)(0.008) + (0.03)(0.992)} = 0.97 \\ P(D=\neg P|R=P) &= 1 - 0.97 = 0.03 \end{aligned}$$

Frequency Table

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

→ $P(\text{Play} | \text{Outlook}) \rightarrow$ what is the prob of play given each outlook

outlook	Play = No	Play = Yes
Sunny	3	2
Overcast	0	4
Rainy	2	3

Bayes Prob

$$\begin{aligned} p(O=S|R=N) &= \frac{3}{5} & p(O=S|R=Y) &= \frac{2}{9} \\ p(O=O|R=N) &= 0 & p(O=O|R=Y) &= \frac{4}{9} \\ p(O=R|R=N) &= \frac{2}{5} & p(O=S|R=Y) &= \frac{3}{9} \\ & 1 & & 1 \end{aligned}$$

$$P(P=Y|e) = \prod_i \frac{P(e_i | P=Y)P(P=Y)}{P(e_i)} \leftarrow P(e_i | P=Y)P(P=Y) + P(e_i | P=\neg Y)P(P=\neg Y)$$

MAP - assign the most probable category h_{MAP} given e_1, \dots, e_n

$$h_{\text{MAP}} = \max_h P(h | e_1, \dots, e_n) = \max_h \frac{P(e_1, \dots, e_n | h)P(h)}{P(e_1, \dots, e_n)} \propto \max_h P(e_1, \dots, e_n | h)P(h)$$

Zero Frequency Problem

Add 1 to the numerator for ever attribute in frequency table so that probability can never be zero

outlook	Play = No	Play = Yes	get add by 1	get add by 3
Sunny	3+1	2+1 → $p(O=S R=N) = \frac{4}{8}$		$p(O=S R=Y) = \frac{5}{12}$
Overcast	0+1	4+1 → $p(O=O R=N) = \frac{1}{8}$		$p(O=O R=Y) = \frac{5}{12}$
Rainy	2+1	3+1 → $p(O=R R=N) = \frac{3}{8}$		$p(O=S R=Y) = \frac{9}{12}$

Missing Values

Missing val is numeric → Sample missing/new data x from $N(\mu_i, \sigma_i)$ for feature i
 Classify the label for that new data $P(y|e_1=x, e_2, e_3, \dots, e_n) \propto P(e_1=x, e_2, \dots, e_n | y) P(y)$

Bayesian Network

Naive Bayes assumption of all feature being indep. to each other is too restrictive

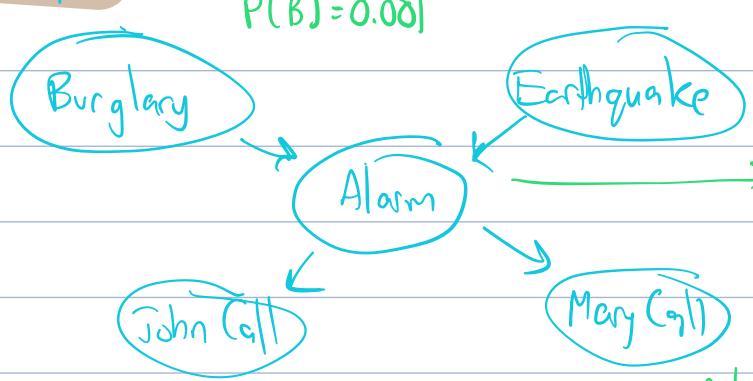
Bayesian Network - conditional indep. among subset of feature

$$P(e_1, \dots, e_n) = \prod_{i=1}^n P(e_i | e_{i-1}, e_{i-2}, \dots, e_1) = \prod_{i=1}^n P(e_i | \text{Parent}(e_i))$$

$$P(A, B, C, D, E) = P(A) P(B|A) P(C|AB) P(D|ABC) P(E|ABC)$$

$P(A|B,C)$ if C is conditionally indep. of A given B , then $P(A|BC) = P(A|B)$

DAG (Graph)



$$P(B) = 0.001$$

$$P(E) = 0.002$$

				$P(A B, E)$
		B	E	$P(A B, E)$
B	E	T	T	0.95
		T	F	0.94
F	T	F	T	0.29
		F	F	0.001

A $P(J A)$	
T	0.9
F	0.05

A $P(M A)$	
T	0.7
F	0.01

$$P(J, M, A, \neg B, \neg E) = P(J|A) P(M|A) P(A|\neg B, \neg E) P(\neg B) P(\neg E)$$

$$= (0.9)(0.7)(0.001)(0.999)(0.998) = 0.00062$$

Lecture 6: Decision Tree

Internal Node - test on an attribute

Branch - subset of the tree from the outcome of the test

Leaf node - doesn't split any further, give class label or distribution

At each node, one attribute is chosen to split training example

Top-down Tree construction

- At start, all training example are at the root.

- Partition the example recursively by choosing one attribute each time

Bottom-up tree Pruning

- Remove subtree or branch from the bottom to top to improve the accuracy on new case
(no overfitting)

Choosing Attribute to split

Information gain - choose attribute that result in greatest info gain

Info gain = info before split - info after split

$$\text{gain}(\text{outlook}) = \text{info}([9,5]) - \text{info}([2,3], [4,0], [3,2])$$

Look at the gain for each attribute → choose attribute w/ the highest gain (purity)

$$\text{Entropy}(p_1, \dots, p_n) = - \sum_i p_i \log_2(p_i) \quad p_i = P(y=c | t) \Rightarrow \text{prob of picking class } c \text{ in data set } t$$

(use ln for the exam instead)

$$\text{gain}(S, a) = \text{entropy}(S) - \sum_{v \in a} \frac{|S_v|}{|S|} \text{entropy}(S_v)$$

↑ # of value of specific type (outlook = Sunny → 5)
↑ subset of data that have specific type of value
and calculate entropy
This is info gain after splitting tree with feature a
↑ All possible type of value for feature a

$$i_E(t) = -P(y=\text{No} | \text{outlook} = \text{Sunny}) \ln P(y=\text{No} | o = \text{Sunny}) - P(y=\text{Yes} | o = \text{Sunny}) \ln P(y=\text{Yes} | o = \text{Sunny}) = -\frac{3}{5} \ln \left(\frac{3}{5} \right) - \frac{2}{5} \ln \left(\frac{2}{5} \right)$$

Splitting for ID Attribute → gain might be high (cos each group is pure) but then there's too many branches

Solution : Gain Ratio reduce bias on high branch attribute

• Takes number and size of branch into account

$$\text{Calculation intrinsic info and gain Ratio}(S, a) = \frac{\text{gain}(S, a)}{\text{intrinsic info}(S, a)}$$

$$\text{info}[1,2,2] = \underbrace{-\frac{1}{5} \ln \left(\frac{1}{5} \right)}_{\text{so 3 branch, 1 data on 1 branch}} - \frac{2}{5} \ln \left(\frac{2}{5} \right) - \frac{2}{5} \ln \left(\frac{2}{5} \right)$$

$$\text{info}[1,2,2] = \text{Entropy} \left(\left(\frac{1}{5}, \frac{2}{5}, \frac{2}{5} \right) \right)$$

so 3 branch, 1 data on 1 branch, 2 data on the other 2 branch

$$\text{info}([2,3], [5,4]) = \frac{5}{14} \text{info}[2,3] + \frac{9}{14} \text{info}[5,4]$$

$$= \frac{5}{14} \left(-\frac{2}{5} \ln \left(\frac{2}{5} \right) - \frac{3}{5} \ln \left(\frac{3}{5} \right) \right) + \frac{9}{14} \left(-\frac{5}{9} \ln \left(\frac{5}{9} \right) - \frac{4}{9} \ln \left(\frac{4}{9} \right) \right)$$

Gini Index

$$\text{Gini}(P) = 1 - \sum_j p_j^2$$

p_j = prob of picking class j

$$\frac{|S_L|}{|S|}$$

same for both side

opt(I)

$$A=400, B=400$$

$$\text{opt(I)} 1 - \left(\frac{300}{400} \right)^2 - \left(\frac{100}{400} \right)^2 = 0.375$$

$$\frac{400}{800}$$

$$(0.375) = 0.1875 \times 2 = 0.375$$

$$A=200, B=200$$

$$A=200, B=200$$

$$\text{opt(II)} 1 - \left(\frac{200}{600} \right)^2 - \left(\frac{400}{600} \right)^2 = 0.44\ldots$$

$$\frac{600}{800}$$

$$(0.44) = 0.33\ldots$$

opt(III)

$$A=400, B=400$$

$$S=800$$

$$\text{right} \rightarrow 1 - \left(\frac{200}{200} \right)^2 - 0^2 = 0$$

$$\frac{200}{800} (0) = 0$$

total data on this split side

$$A=200, B=200$$

$$S_V=600$$

$$A=200, B=0$$

$$S_V=200$$

$$\Delta_i(t) = 0.5 - 0.375 \rightarrow 0.125$$

$$\Delta_i(t) = 0.5 - 0.33 = 0.17 \rightarrow \text{more gain!}$$

DT w/ Numeric Attribute

Numerical Attr. can be use several time in DT while Nominal (class) Attr. only once

Place split points halfway between the value eg. 71 ; 72 ...

No : No

split here < 71.5 and ≥ 71.5

$$p(\text{No} < 71.5) \quad p(\text{No} \geq 71.5)$$

Handle missing value in training

1. Ignore data w/ missing value

2. Ignore attribute ~~it~~

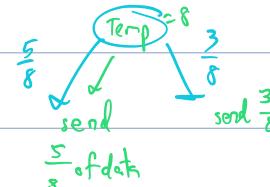
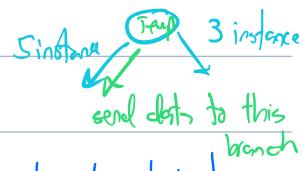
3. Treat missing value as another nominal value (eg. label as "error")

4. Estimated missing value

Classifying data w/ missing value in the attr.

1. follow the leader \rightarrow send all the data w/ missing value on the attribute to the branch w/ most instance

2. "Partition" the instance \rightarrow Send down proportion of data to each branch based on the proportion of training instance



What can you do to make DT smaller and easier to interpret?

1. Pre Pruning (Early stop)
2. Post pruning
3. Limit depth of DT
4. use gain ratio (could help)

Lecture 7 : Data Preparation and Causal Inference

reduce # at branch)

In DT : Prepruning and Post pruning

Prepruning - stop splitting the tree if $g(s, a) < \text{threshold}$

Postpruning - Construct a complete DT first w/o threshold, then prune it back

- When prune \rightarrow replace subtree w/ a single leaf node (remove all children of node but not the node itself)

Prune the node if validation error rate is reduce after pruning $\text{err}(\text{before prune}) > \text{err}(\text{after prune})$

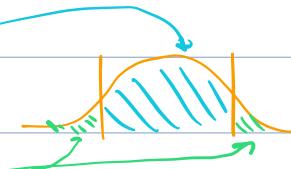
(hold out)
The observed error rate $f = \frac{\text{errors}}{n}$ follows Normal distribution

Central limit theorem - the standardised avg of any population of i.i.d RV X_i w/ μ_x and $\sigma^2 \sim N(0, 1)$

Confidence Interval of Normal Dist.

% confidence interval for RV X w/ $\mu=0$ = $P(-z \leq X \leq z) = x$

and $P(-z \leq X \leq z) = 1 - 2P(z \geq |X|)$
tail distri.



Data Preparation

Data cleaning

1. Join relational data file together
2. Deal w/ missing entry \rightarrow ignore data w/ missing entry
 - \rightarrow Treat missing entry as separate value
 - \rightarrow fill missy entry w/ mean or median value

Conversion

Convert Ordinal data (ordered attri (grade, size)) to numeric preserving natural order \rightarrow ordinal encoding

Convert nominal data (unordered) to binary dummy variable with one-hot encoding \rightarrow Dummy variable encoding

Discretization (Binning) - reduce number of values for a continuous attribute

\rightarrow Convert Numerical attr. to discrete attr. \rightarrow group them into bucket

4. Build balance train & val sets (no imbalance set)

5. Subset Selection - Select subset of attribute that is the most relevant

- Backward elimination - Use everything then remove attr. one by one
- Forward selection - Choose one attr. at a time

Causal Inference

2 attr. are related but did one attr. have the effect on the other? (Correlation doesn't imply Causation)

Sample Selection bias → data collection process is bias

1) Experimental studies - Researcher intervene and observe what happens, indep. var are controlled

1.1) Randomized Controlled Trials (RCT) → (A/B Testing)

Randomized exp. where each subject is randomly assigned to a treated group or a control group

in order to control for extraneous factors. Randomization mitigate the sample selection bias

Extraneous factor - variable that you're not investigating (not indep. var) that can potentially affect the result

1.2) Quasi-Experiments

Compare groups and measure effects w/o randomization of the subjects.

Assignment of subjects is not random → Selection bias is an issue.

2) Observational Studies - Researcher studies what occurs but doesn't change the subject

- Indep. var is not under control of the research (opposite of Experimental study)

- Selection bias is a concern

2.1) Cross-Sectional Study

- Involve data collection from a population or subset of it at one specific point of time

e.g. Sample of commuter on a given morning and study the mode of transportation

2.2) Panel (or cohort) studies

Observe a group of subjects over a span of time.

2.3) Case-control study

Study in which 2 existing groups differing in outcome are identified and compared on the basis of some supposed causal attribute

e.g. Compare history of cancer patient vs. non-patient.

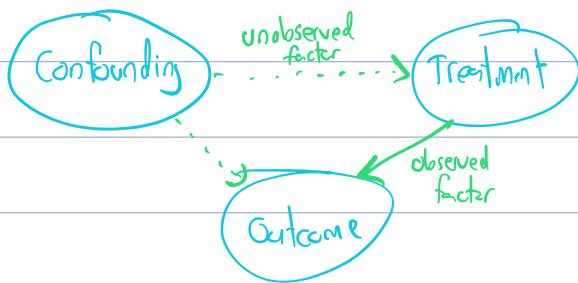
We want to identify causal, but there might be confounding variable

Confounding variable

is an extraneous var. that correlate (directly or inversely) w/ both dep and indep. var.

in a way that "explain away" some or all of the correlation between the dep. and indep. var

- Confounding is part of exogeneity (one of assumption in Gauss-Markov)



How to deal w/ confounding var. (to identify causal effects)

1. RCT

Lab experiment → lead to low external validity (cannot explain data outside of observed population)

Create situation w/ desired conditions

Manipulate some var. while controlling others

Examine dependent variable

Field experiment → lead to high external validity

Research study in natural setting, manipulate some variable (no control), examine dep. var.

Randomized Exp. (lead to high internal)

Quasi-Exp. (lead to low internal) more confounding var. from selection bias

Field

high internal / high external validity

low internal / high external validity

Lab

high internal / low external validity

low internal / low external validity

2. Fixed Effects Models → for panel study

Assume confounding var. has a fixed or constant relationship w/ the dep. variable across all observation

(So it can be subtracted out w/ other data from diff group)

3. Difference in difference models (quasi-exp.)

We take a common trend assumption (treatment & control group has the same overall trend)

e.g. Find causal effect on new law being introduced

City w/o new law (Control), Cities w/ new law (Treatment)

Year	Treatment	Control
2019	9	10.2
2021	6.9	8.7

We use control group

to find the trend of the data

$$(Y_{T_2} - Y_{T_1}) - (Y_{C_2} - Y_{C_1}) = (6.9 - 9) - (8.7 - 10.2) = -0.6 \rightarrow \text{decrease after law}$$

4. Propensity Score Matching (PSM) (Cross Selectional data)

Lecture 8 : Evaluation

Bias-Variance Trade off

Bias - Error from inaccurate prediction b/c it's too simple (not enough param)

Variance - error from sensitive to small fluctuation. To many params.

Under fitting - model is too simple to represent all relevant characteristics

Over fitting - model is too complex and fit irrelevant/noise characteristics

Resampling Method

• Holdout procedure

- reserve data for test set and use remaining for training and validation

• Stratified holdout

- guarantee that classes are proportionally represented in the test/train set

• Repeated holdout

Randomly select test set (not val set like cross k-fold) and avg the error rate estimate

K-fold Cross-Validation

- fixed # of k partition of the data (fold)

- Each fold, select one partition as a validation set and the rest as training set (each partition is used for validation only once)

- Error rate is estimate by avg of error rate of all k-fold

Use paired t-test to check if estimate error of model is "significantly" difference from esti. error rate of another model

Bootstrap

1. Draw a random sample of size n and calculate statistic (error rate, mean) with these data

2. Repeat it k many times, getting $\text{error} = (e_1, e_2, e_3, \dots, e_k)$ $\mu = (\mu_1, \mu_2, \mu_3, \dots, \mu_k)$

3. Calculate variance of your statistic to learn about population statistics

4. Calculate confidence interval of your original data

Measuring Error

		Predicted	
		T	F
Actual	T	TP	FN
	F	FP	TN

$$\text{Error rate} = \frac{\# \text{ of error}}{\# \text{ of instance}} = \frac{(FN+FP)}{N}$$

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

$$\text{Recall} = \frac{TP}{TP+FN}$$

$$\text{Precision} = \frac{TP}{TP+FP}$$

$$\begin{aligned}\text{Specificity} &= \frac{TN}{TN+FP} \\ \text{False Alarm} &= \frac{FP}{FP+TN}\end{aligned}$$

Gain Curve

% of population vs # of target (that you're interested in)

$$\text{gain curve} = \frac{\# \text{ of target for } \% \text{ of population}}{\text{Total } \# \text{ of target}} = \frac{n}{N}$$

Lift Curve

$$\text{Lift}(q) = \frac{\text{Gain}(q)}{q}$$

$q = \% \text{ of sample size}$

ROC Curve

TP rate (Recall) vs False Positive



(look at exercise on how to calculate TPR and FPR)

With gain curve & lift curve, you can see which model perform best at low sample size
(like it predict 80% correct w/ 50% data vs predict 80% at 80% data)

Lecture 9 : Ensemble Methods

Idea - Combining Multiple Models

Advantage

- Often Improve Predictive Perf.

Disadvantage

- Output is hard to analyze

Bagging

- Aim to reduce variance
- Each Model receive equal weight

- How it work :

- 1) Sample several training set of size n from the training data (Bootstrap)
- 2) Built a classifier (model) for each training set
- 3) Combine classifier's prediction

- Advantage:

- Can be apply to numeric prediction & classification
- Help perf if data is noisy
- The more model, the better perf (w/ diminishing return)
- Can be parallelized cus model can be create independently

Random Forest → Combine Bagging with Decision Tree

→ Randomized training data (Bootstrap) and feature!

(randomize the feature selection for each DT for bagging)

For each data x_i , return class that was predicted most often for all DT.

Each tree has high variance & low bias → Avg all the tree help reduce variance

Boosting

Idea: Combine several weak learners into a strong learner

New models are influenced by perf. of previously built model (model try to fix the problem of previous model)

- Increase capability of the model, minimizing the bias (training error)
- Weights models according to performance
eg. Ada Boost, XG Boost

Bias-Variance Tradeoff!

Bagging reduces variance of low-bias models

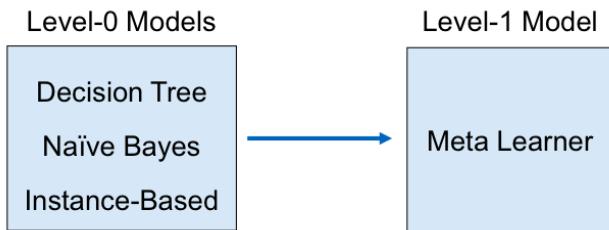
- Low-bias models are "complex" and unstable
- Bagging averages them together to create stability

Boosting reduces bias of low-variance models

- Low-variance models are simple with high bias
- Boosting trains sequence of models on residual error → sum of simple models is accurate

Stacking

Prediction from each homogenous classifiers are used as an input into a meta learner which attempt to combine prediction to create a final best prediction



Clustering - Unsupervised learning

- Item in same cluster are similar (high similarity)
- Item from diff cluster are different (low similarity)

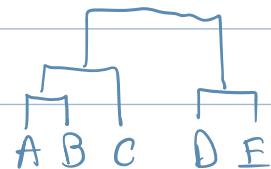
1. Hierarchical Clustering

Bottom Up

Start from each data has its own cluster

At each step, merge 2 closest clusters

use Dendrogram



Top down

Start w/ all data in a cluster

Separate into 2 cluster at a time

Use Dissimilarity Matrices (calculate based on distance) to find closest cluster

MST - Minimum Spanning Tree

- Cheapest path from data A to B → A and B is in same cluster
- Each cheapest path between 2 data/cluster is a join in the cluster.
- Create a dendrogram from the cluster

K-Means Clustering

- work w/ numeric only

How it works:

- 1) initialise k cluster centroid in random
- 2) Assign every data to the nearest cluster (Euclidean- L_2 , Manhattan- L_1 , etc.)
- 3) Move the centroid to the mean of its assigned data
- 4) Repeat 2 & 3.

Result can vary significantly depending on initialised centroid

Probability-Based Clustering (Gaussian Mixture Model)

Mixture Model = Model each cluster w/ prob. distribution (Gaussian)

Each probability distribution give prob of the data belonging to each cluster $p(x_i) = f(c_1, c_2, \dots, c_K)$

Expectation Maximization (algorithm to optimize GMM)

Start w/ initial param $\sigma_A, \sigma_B, \mu_A, \mu_B, p(A)$

Iterate between E and M

Expectation - calculate probability of data being on each cluster $\delta_t(z_{ik}) = p(z_{ik} | x_i, \pi_t, \mu_k, \sigma_k)$

Maximization - Find the maximum likelihood of each param from $\delta_t(z_{ik}) - \pi_{k+1}^*, \mu_{k+1}^*, \sigma_{k+1}^*$

Use the recalculated param to calculate expectation again $\delta_{t+1}(z_{ik})$

Use recalculate expectation $\delta_{t+1}(z_{ik})$ to find MLE param again.

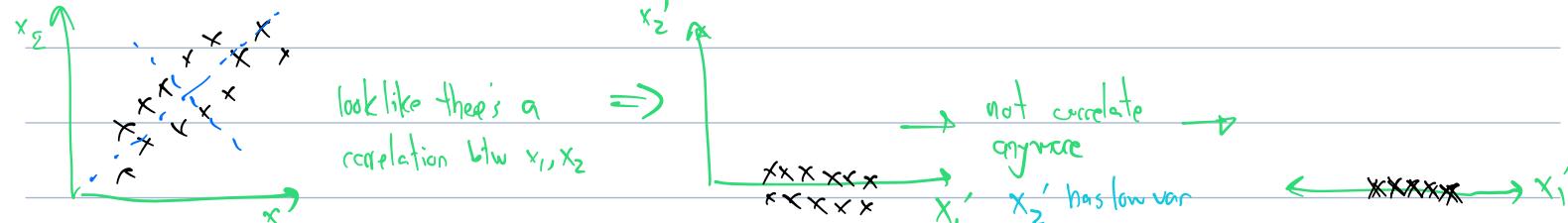
Lecture 10: Dimensionality Reduction

PCA - convert a set of possibly correlated variables into a set of values of linearly uncorrelated var.

This can help combat multicollinearity in the regression analysis

Transform the coordinate system to the principle component coordinate system

eigenvector \vec{v} are direction that doesn't rotate when transform only scaling b \Rightarrow



How to transform data w/ PCA

, can get rid

Step 1: Compute zero mean data (compute mean for each column), then $\tilde{X} = X - \bar{X}$

Step 2: Calculate Covariance Matrix of \tilde{X}

$$\Sigma_x = \frac{\tilde{X}^T \tilde{X}}{N-1} \quad \text{or} \quad \Sigma_x = \begin{pmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) \end{pmatrix}$$

$$\text{Var}(x_i) = \frac{1}{N-1} \sum_j (x_{ji} - \bar{x}_i)^2 \quad \text{if } x \text{ is zero mean}$$

$$\text{Var}(x_i) = \frac{1}{N-1} \sum_j x_{ji}^2$$

$$\text{Cov}(x_i, x_j) = \frac{1}{N-1} \sum_k (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j) = \frac{1}{N-1} \sum_k x_{ki} x_{kj}$$

Step 3: Calculate ν and σ of $\Sigma_x \rightarrow \Sigma = \Phi \Lambda \Phi^T$

Final data: $Z = X \cdot \Phi$ where Φ = eigenvector matrix

Can leave out dimension w/ small variance \rightarrow have enough σ to cover 80-90% variance

Reconstruct data $X = Z \cdot \Phi^T + \text{original mean}$

SVD

$$A = U \Sigma V^T \quad \text{where} \quad \Sigma(A) = \sqrt{\lambda(\Lambda A \Lambda)}$$

V is the principle component coordinate system

$U \Sigma$ = coordinate of A in the new coordinate system • scaling

$$Z = A V = U \Sigma$$

PCA is not feature selection! it's feature extraction.

Feature selection, you reduce # of dim by throwing away irrelevant information (use same data but throw away some feature)

Feature Extraction (PCA), reduce # of dim by transforming data (projecting) to low dim data, thus still maintaining the same information

Lecture 11 : Convex Optimization & NN

$$\nabla f(x) = \left(\begin{array}{c} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{array} \right) \quad \text{- gradient of } f = \text{direction of steepest ascent}$$

$$-\nabla f(x) = \text{steepest descent}$$

Convex

f is convex iff $\forall \gamma \in [0, 1] \quad f(\gamma x + (1-\gamma)y) \leq \gamma f(x) + (1-\gamma)f(y)$

$f(y) \leq f(x) + \nabla f(x)^T(y-x) \quad \forall x, y \in \text{dom}(f), \text{ dom}(f) \text{ is convex}$

f is convex iff $\nabla^2 f(x) \geq 0$ or $x^T H(x, y, \dots) x \geq 0$ (Hessian is PSD)

f is convex $\Leftrightarrow H_f(x)$ is PSD

if $\nabla f(x) = 0, H_f(x)$ is Positive Definite $\rightarrow x$ is a local min & $\nabla f(x) \geq 0$

$\nabla f(x) = 0, H_f(x)$ is Negative Definite $\rightarrow x$ is a local max & $\nabla f(x) \leq 0$

$\nabla f(x) = 0, H_f(x)$ is indefinite $\rightarrow \nabla f(x)$ has both < 0 and $> 0 \rightarrow$ saddle point

$\nabla f(x) = 0, H_f(x)$ is PSD \rightarrow inconclusive $\nabla f(x) \geq 0$

but

x is local min $\rightarrow H_f(x)$ is PSD

x is local max $\rightarrow H_f(x)$ is PSD

Gradient Descent

$$x^{(k+1)} = x^{(k)} - \alpha \nabla f(x^{(k)}) \quad \text{change over iteration}$$

α can be a line search (expensive) or fixed α or α_t

Start at $x^{(1)}$ \rightarrow calculate $\nabla f(x^{(1)})$ (min at $x^{(1)}$ by $\nabla f(x^{(1)}) = 0$)

$$\text{then } x^{(2)} = x^{(1)} - \alpha \nabla f(x^{(1)})$$

$$\text{if } \alpha \text{ not given, then } \alpha^* = \arg \min_{\alpha} f(x^{(1)} - \alpha \nabla f(x^{(1)})) = \frac{\partial}{\partial \alpha} f(x^{(1)} - \alpha \nabla f(x^{(1)})) =$$

input for func.

Momentum Gradient Descent

$$\text{Use momentum } x^{(k+1)} = x^{(k)} - \alpha V_t \quad \text{and } V_t = \beta \nabla f(x^{(k)}) + (1-\beta) V_{t-1}$$

Increase the step if the direction go in the same direction (accumulate gradient)

Multiple Linear Regression Model

+1 for bias

Estimated equation

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_d x_d = \underbrace{x \hat{\beta}}_{\text{error}} + e$$

$y \in \mathbb{R}^{n \times 1}$ $x \in \mathbb{R}^{n \times (d+1)}$ $\hat{\beta} \in \mathbb{R}^{(d+1) \times 1}$ $e \in \mathbb{R}^{n \times 1}$

Loss function of MLR = RSS

$$RSS = e^T e = (y - x \hat{\beta})^T (y - x \hat{\beta}) \rightarrow \text{convex b/c let } V = y - x \hat{\beta} \quad V^T V = (Vx)^T (Vx) = \|Vx\|_2^2 \geq 0$$

Ordinary Least square estimation → find param that will minimize the loss

$$\beta_{LS} = (x^T x)^{-1} x^T y \quad \text{so} \quad \hat{y} = x (x^T x)^{-1} x^T y$$

Do OLS Estimation via GD

$$\text{Use mean square error func} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

$$\frac{\partial L}{\partial \beta_0} = -\frac{2}{n} \sum_i (y_i - \hat{y}_i) \quad \frac{\partial L}{\partial \beta_i} = -\frac{2}{n} \sum_i x_i (y_i - \hat{y}_i)$$

$$\hat{\beta}_0 = \hat{\beta}_0 - \alpha \frac{\partial L}{\partial \hat{\beta}_0} \quad \hat{\beta}_i = \hat{\beta}_i - \alpha \frac{\partial L}{\partial \hat{\beta}_i}$$

Best Linear unbiased estimator

Recap Gauss-Markov Theorem : OLS estimator give the BLUE of the coefficient

Unbiased = $E(\hat{\beta}_j) = \beta_j$ Best = give lowest variance compare to other linear unbiased estimator

But often, there exist biased estimator that give lower error/variance

Regularization — estimator will be more biased but lower var

Ridge Regression (L_2)

$$L(\beta) = \sum_i (y_i - \hat{y}_i)^2 + \lambda \sum_j \beta_j^2 \quad \beta \text{ will shrink all equally}$$

Lasso (L_1)

$$L(\beta) = \sum_i (y_i - \hat{y}_i)^2 + \lambda \sum_j |\beta_j| \quad \beta \text{ will be more sparse (more 0's entry)}$$

b/c less important feature get set to 0, Lasso

perform a feature selection. Also, make model

simpler to interpret.

Logistic Regression

$$\text{Likelihood } L = \prod_i p_i^{y_i} (1-p_i)^{1-y_i} \quad \text{let } p(y|x) = \sigma(\beta_0 + \beta_1 x) = \sigma(z)$$

$$\text{Log Likelihood } LL = \sum_i y_i \ln p_i + (1-y_i) \ln (1-p_i)$$

$$\text{Negative LL} = -\sum_i y_i \ln p_i - (1-y_i) \ln (1-p_i)$$

Is NLL convex?

$\beta_0 + \beta_1 x$ is affine func \rightarrow convex

$\sigma(\beta_0 + \beta_1 x)$ is convex

is $f(g(z)) = -\ln \sigma(\beta_0 + \beta_1 x)$ convex?

$$\ln(1 + \exp(-(\beta_0 + \beta_1 x)))$$

$$g''(z) > 0$$

* if $f''(g(z)) \geq 0$ (f is convex) and $f'(g(z))$ is monotonically increasing and g is convex
then $f(g(z))$ is convex

\therefore You can find the minimum of NLL

Prove

$$\begin{aligned} f'(g(z)) \cdot g''(z) &= f'(g(z)) \cdot g''(z) + g'(z) f''(g(z)) \cdot g'(z) \\ &= f'(g(z)) \cdot g''(z) + g'(z) f''(g(z)) \geq 0 \end{aligned}$$

Lecture 12: NN

Risk function R is a expectation of a loss function (avg loss per data point)

Loss func - loss of the whole dataset $\rightarrow L(\theta) = -\sum_{i=1}^n p_i \log(\hat{y}_i) + (1-p_i) \log(1-\hat{y}_i)$

$$R(\theta) = \frac{1}{N} \sum_n L(y_n, g(\theta, x_n))$$

if Loss func = MSE,

$$\text{then } R(\theta) = \frac{1}{N} \sum_n \left(\frac{1}{2} (y_n - g(\theta, x_n))^2 \right)$$

$$L(\theta) = -\frac{1}{n} \sum_{i=1}^n p_i \log(\hat{y}_i) + (1-p_i) \log(1-\hat{y}_i)$$

* We use Risk func (avg w/ $\frac{1}{n}$)

here in this class

GD

Update w/ small $\alpha \rightarrow$ very small step, might get stuck in saddle point

w/ large $\alpha \rightarrow$ oscillate

\rightarrow In ML

we use bkt loss

A 1 layer NN is similar to Logistic Regression. But NN can stack many layers = more complex

no hidden layer

output layer = 1 layer

Input layer = # of feature of data

Output layer = # of class for layer for layer

let model consist of H_1 = hidden node¹, H_2 = hidden node², Y = output node, X = input node

of weight in node : $(X+1) \times H_1 + (H_1+1) \times H_2 + (H_2+1) \times Y$

Backpropagation

Calculate gradient of params of NN w/ respect to Loss func

$$w_{ij}^{(l+1)} = w_{ij}^l - \alpha \frac{\partial R}{\partial w_{ij}} \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad W^{(l+1)} = W^{(l)} - \alpha \nabla_W L(W^{(l)})$$

$$w_{jk}^{(l+1)} = w_{jk}^l - \alpha \frac{\partial R}{\partial w_{jk}}$$

Matrix Notation for 2 layer NN (1 hidden layer and σ activation func)

$$\left. \begin{array}{l} z^{(1)} = XW^{(1)} + b^{(1)} \\ A = \sigma(z^{(1)}) \\ z^{(2)} = AW^{(2)} + b^{(2)} \\ \hat{y} = \sigma(z^{(2)}) \end{array} \right\} \quad \hat{y} = \sigma(\sigma(XW^{(1)} + b^{(1)})W^{(2)} + b^{(2)})$$

Mini Batch

with n data, create m batch, each w/ $\frac{n}{m}$ data. For each batch, do forward pass, calculate loss, backward pass and update param. Update param m time per epoch.

Stochastic Gradient

Minibatch w/ only one data for each batch, calculate loss of 1 sample and do backprop and update the parameter. So if $n=500$, do 500 update per epoch.

Both provide noisy gradient which can be good for NN.

Hidden layer do feature extraction (transform info, didn't get rid of info (feature selection))

Universal Approximation Theorem

NN w/ one hidden layer and non linear activation func can approx. any func.

Lecture 13: Reinforcement Learning

What is RL?

RL incorporate time into learning.

It learn to make sequential decision in an environment which maximize "rewards"

- Reward tell how well our agent behaved (reward can be positive or negative)

How RL diff from Supervised and Unsupervised L?

- No supervision
- Feedback can be delayed
- Time matter (sequential data)
- Not i.i.d. (Action affect subsequent data)

Markov Process - Observe sequence of states or a chain

System change between states according to some dynamics

Markov Process iff future system dynamic from any state depend only on current state.

"Future is independent of the past, given the present")

Markov Decision Process - Markov Process w/ action and reward

(if only 1 action exist for each state and all reward is same, then it becomes Markov Process)