

Classification

- uncertainty:
 - deterministic: when $X > Y, Y > Z \rightarrow X > Z$
 - but in probabilistic: we consider 4 dices:

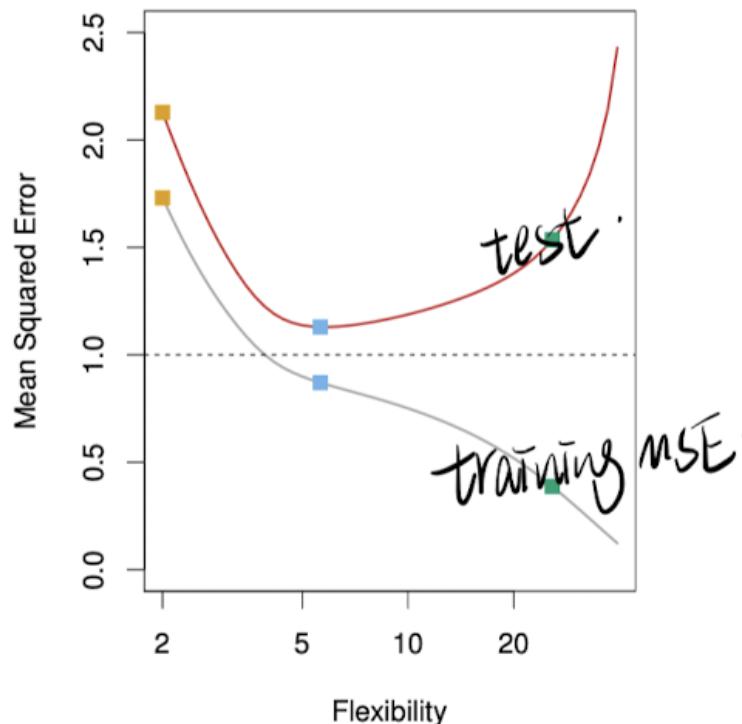
A: 0 0 4 4 4 4
B: 3
C: 2 2 2 2 6 6 $C > D : \frac{2}{3} \times$
D: 1 1 1 5 5 5
 $P(A > B) : 2/3$ $D > A : \frac{1}{2}$
 $P(B > C) : 2/3$
 $P(C > D) : 2/3$
 $P(D > A) : 2/3 - \text{no transitivity}$

- Two methods:
 - Bayes inference: only know the past data and distribution ,and we can get new information, then we can update the distribution, e.g. weather
 - Frequency: on the same condition and repeat many many times
 - pros: easy to understand
 - cons: only suitable for some special cases
- Statistical machine learning
 - $Y = f(X) + \epsilon$
 - reasons:
 - prediction
 - error
 - reducible:
 - irreducible: Y is also a function of ϵ , this can't be predicted by X
 - inference
 - estimate f
 - Parametric
 - procedure
 - assum the functional form
 - estimate parameters
 - cons: assumption may not correct
 - how to solve this? choose flexible models<- but more complex model can lead to overfitting

- Non-parametric: 不是没有参数，而是没有预先设定的参数形式
 - cons: require a lot of data(far more than parameter model)
- Assessing model accuracy
 - measure the quality of fit
 - MSE: most test MSE(注意没有irreducible error)

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2,$$

- plot



- when training MSE decreases, testing MSE increases --> overfitting<--think some random things as the reason.
- smoothing spline

$$\sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

- $\lambda \rightarrow \infty$, 不能有太多的二阶导的影响, 所以是linear function
- $\lambda \rightarrow 0$, more flexible function
- 根据曲线balance between underfitting and overfitting

- the smoothing spline seems to get min testing MSE, but why we still need some other methods? In principle, it's perfect, but in practice, it's not perfect especially in the high dimensional case.
- The Bias-Variance Trade-Off

- formula

$$\min_{\hat{f}} \left[f(x_0) - \hat{f}(x_0; \mathcal{D}) \middle| X = x_0 \right]^2 ?$$

- \hat{f} 是根据 training data 训练出来的, x_0 不是 training data 里的数据 (因为我们想要的是 min test MSE), 且我们不需要一个只适用于某一个特定 dataset 的模型, 而应该是最好所有的 dataset (training set), 在这个想法基础上, 我们需要加上期望
- MSE

$$\hat{g} = \arg \min [\sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g'(t)^2 dt]$$

training error. $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{g}(x_i))^2$

empirical testing error $\frac{1}{n} \sum_{i \in D_{test}} (y_i - \hat{g}(x_i))^2$

expected testing error. $E[(Y - \hat{g}(X)) | D_{train}]^2$

when $n_{test} \rightarrow \infty$, empirical testing error \rightarrow
expected testing error.

- 事实上, 有没有条件是有一定区别的, 但非常非常小, 可以忽略, 因此我们不考虑是否存在 condition
- 化简

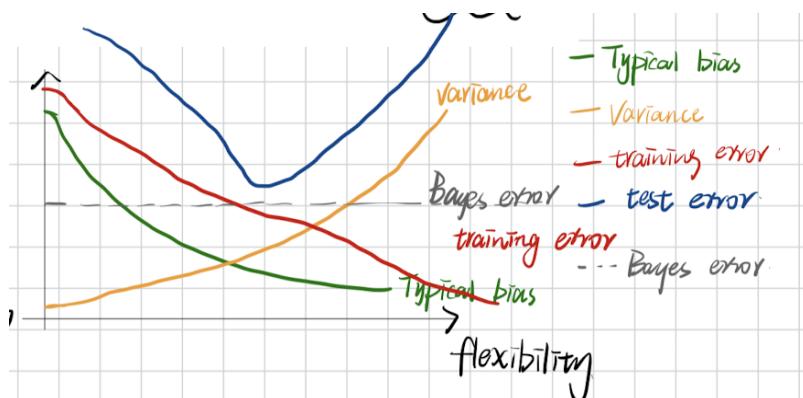
$$\begin{aligned} & E[f(x_0) - \hat{f}(x_0; \mathcal{D}) | X = x_0]^2 \\ &= [f(x_0) - E_D[\hat{f}(x_0; \mathcal{D})] + E_D[\hat{f}(x_0; \mathcal{D})] - \hat{f}(x_0; \mathcal{D}) | X = x_0]^2 \\ &= [f(x_0) - E_D[\hat{f}(x_0; \mathcal{D})] | X = x_0]^2 \\ &\quad + [E_D[\hat{f}(x_0; \mathcal{D})] - \hat{f}(x_0; \mathcal{D}) | X = x_0]^2 \\ &\quad + 2(f(x_0) - E_D[\hat{f}(x_0; \mathcal{D})]) [X = x_0] (E_D[\hat{f}(x_0; \mathcal{D})] - \hat{f}(x_0; \mathcal{D})) | X = x_0 \\ &\text{when we take expectation. } E_D(E_D[\hat{f}(x_0; \mathcal{D})]) = E[\hat{f}(x_0; \mathcal{D})] = 0. \\ & E[f(x_0) - \hat{f}(x_0; \mathcal{D}) | X = x_0]^2 = E_D[f(x_0) - E_D[\hat{f}(x_0; \mathcal{D})] | X = x_0]^2 \text{ 在 } X=x_0 \text{ 情况下, } f(x_0) \text{ is fixed.} \\ & \text{Variance: } E_D[\hat{f}(x_0; \mathcal{D})] \text{ is fixed} \Rightarrow \text{delete } E_D[\hat{f}(x_0; \mathcal{D})]. \\ & + [E_D[\hat{f}(x_0; \mathcal{D})] - \hat{f}(x_0; \mathcal{D})]^2. \quad \hat{f} \text{ is random, because sample is random.} \\ &= [\hat{f}(x_0) - E_D[\hat{f}(x_0; \mathcal{D})] | X = x_0]^2. \quad \hat{f} \Rightarrow \text{对固定 } x_0 \Rightarrow \text{对不同 } D \text{ 的 } \hat{f}(x_0; \mathcal{D}). \quad (\text{因为 } \hat{f} \text{ 不同的 } D \text{ 不同).} \\ &+ E_D[\hat{f}(x_0; \mathcal{D})] \hat{f}(x_0; \mathcal{D}) | X = x_0 \quad \text{样本内值, 得到 } E_D[\hat{f}(x_0; \mathcal{D})]. \quad \text{不断由抽样取 } D \text{ 的 } \hat{f}(x_0; \mathcal{D}). \\ &= \text{Bias}^2 + \text{Variance}. \end{aligned}$$

- 思考

在作业中，train set 是固定的，true y 不是。
 $\therefore \text{train set} \begin{cases} y_{\text{set } 1} \rightarrow \hat{f}_1 \\ y_{\text{set } 2} \rightarrow \hat{f}_2 \\ y_{\text{set } 3} \rightarrow \hat{f}_3 \\ \vdots \\ y_{\text{set } 100} \rightarrow \hat{f}_{100} \end{cases}$

对每个 $x_0 \in D_{\text{test}}$, $\hat{f}(x_0) = \frac{\hat{f}_1(x_0) + \dots + \hat{f}_{100}(x_0)}{100} = E_D[\hat{f}(x_0|D)]$,
 $\Rightarrow \text{Var}(\hat{f}(x_0|D))$.
 此时 x_0 也是固定的，每个 x_0 有一个 Var .
 再求 E , 针对 x_0 求均值 \rightarrow 平均的 Variance .
 f_{true} 也是，~~但是~~ 对 100 方差，再对 500 均值.

- 这里的bias指的是，在固定 x_0 的情况下，我用的模型和真实值差多少
- Variance指的是，在固定 x_0 的情况下，当我用同一个大的dataset，但是不同的training set的情况下， \hat{f} 也就会变，他的变化浮动大不大，也就是对不同的training set的敏感程度如何，
- 例如：a constant model, $f=c$, 不同得training set得到的model相同，variance=0，但bias会非常大；a very flexible model, 对training set的敏感程度很大，variance很大，bias会比较小--因为比较准确
- plot



- Why machine learning method is flexible but small error?
 - 1. signal + noise, signal is very large than the noise--> unlikely to overfitting
 - K line?
- 不是说如果true is non-linear 就不能用linear model to estimate, 如果bias最小，那么也是可以用的
- linear model
 - simple linear regression
 - concept

population level: $Y = \beta_0 + \beta_1 X + \varepsilon$.

↓ sampling frequent perspective.

sample level: $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$.

$$\langle \hat{\beta}_0, \hat{\beta}_1 \rangle = \arg \min \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2.$$

$(\hat{\beta}_0, \hat{\beta}_1)$ is the realization of (β_0, β_1)

$\langle \hat{\beta}_0, \hat{\beta}_1 \rangle$ is random variables depending on sample.

- LSE: least squares criterion

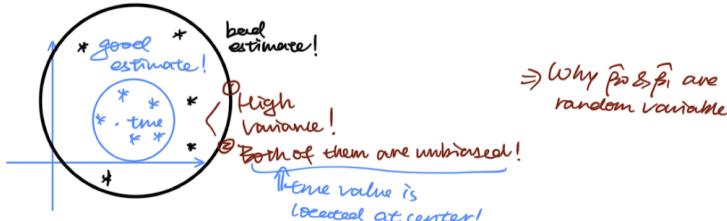
- $RSS = \sum e_i^2$
- coefficient estimate

► The least squares coefficient estimates are given as

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

where $\bar{y} \equiv \frac{1}{n} \sum_i y_i$ and $\bar{x} \equiv \sum_i x_i$.

- Assess the accuracy



- SE: 标准误, 一般等于 $\frac{\sigma^2}{n}$, 在这里, 为什么减2, 因为是2个coefficient, 参见LCM

$$SE(\hat{\beta}_0)^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right], SE(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad (7)$$

where σ^2 can be estimated by $\hat{\sigma}^2 = \frac{RSS}{n-2}$.

Accordingly, we have 95% confidence interval for $\hat{\beta}_0, \hat{\beta}_1$ estimators $\hat{\beta}_0 \pm 2se, \hat{\beta}_1 \pm 2se$.

- 注: 我们可以说 $P(\hat{\beta} - 2se \leq \beta \leq \hat{\beta} + 2se) = 0.95$, 但如果我们一旦知道了 $\beta = 1.1, se = 0.2 \rightarrow P(0.7 \leq \beta \leq 1.5) = 0.95$ 这是不对的, 因为此时, β 内的数字都是固定的, 不是random variable, 意思是说, 我们有95%的把握认为真实参数是在这个区间内, 而不是真实参数在这个区间的概率是95%

- $RSE = \sqrt{\frac{1}{n-2} RSS}$ residual standard error

- $R^2 = \frac{TSS - RSS}{TSS}, TSS = \sum (y_i - \bar{y})^2$

- for training data, non negative
- > 1 better
- only useful for simple linear regression
- just focus on training error

- for testing data it may be negative -- overfitting
- multiple linear regression
 - concept: $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$
 - estimate
 - method 1

$$\hat{\beta} = \arg \min_{\beta} (\vec{y} - \vec{x}\beta)^T (\vec{y} - \vec{x}\beta)$$

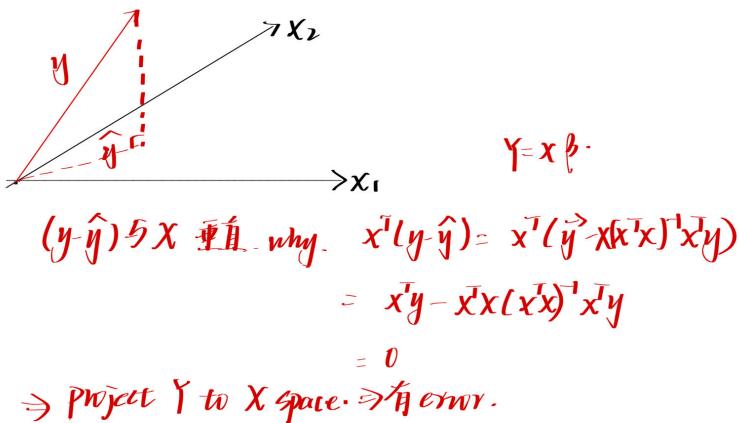
$$\frac{\partial (\vec{y} - \vec{x}\beta)^T (\vec{y} - \vec{x}\beta)}{\partial \beta} = \frac{\partial (\vec{y} - \vec{x}\beta)}{\partial \beta} \frac{\partial (\vec{y} - \vec{x}\beta)^T (\vec{y} - \vec{x}\beta)}{\partial (\vec{y} - \vec{x}\beta)} \vec{x}$$

$$= -\vec{x}^T (\vec{y} - \vec{x}\beta) = 0$$

$$\vec{x}^T \vec{y} - \vec{x}^T \vec{x}\beta = 0$$

$$\beta = (\vec{x}^T \vec{x})^{-1} (\vec{x}^T \vec{y})$$

- method 2
- understanding using projection



- 在assumption下, unbiased and $Var(\hat{\beta})$ 固定

$\vec{Y} = \beta^T \vec{x} + \varepsilon$
not design matrix

Assumption ① $\vec{Y} = \beta^T \vec{x} + \varepsilon$ linear
② ε independent noise $E[\varepsilon_i \varepsilon_j] = 0$

$E[\hat{\beta}] = E[(\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{Y}]$

design matrix (collection of sample values)

$= E[(\vec{x}^T \vec{x})^{-1} \vec{x}^T (\vec{x}\beta + \varepsilon)]$

$= \vec{\beta}$ unbiased!

\vec{Y} simple model is large bias (because this simple one could be non-linear).

because you put intercept into \vec{x}

$Var(\hat{\beta}) = Var(\beta + (\vec{x}^T \vec{x})^{-1} \vec{x}^T \varepsilon)$

$$= Var((\vec{x}^T \vec{x})^{-1} \vec{x}^T \varepsilon)$$

$$= E[(\vec{x}^T \vec{x})^{-1} \vec{x}^T \varepsilon \varepsilon^T (\vec{x}^T \vec{x})^{-1}]$$

$$= E[(\vec{x}^T \vec{x})^{-1} \vec{x}^T \varepsilon]^T E[(\vec{x}^T \vec{x})^{-1} \varepsilon]$$

$$= 0$$

Assume ③ Constant Variance σ^2

$$= (\vec{x}^T \vec{x})^{-1} \vec{x}^T E[\varepsilon \varepsilon^T] \vec{x} (\vec{x}^T \vec{x})^{-1}$$

$$= (\vec{x}^T \vec{x})^{-1} \vec{x}^T (\vec{x}^T \vec{x})^{-1} E[\varepsilon \varepsilon^T]$$

$$= (\vec{x}^T \vec{x})^{-1} \sigma^2 = \sigma^2 (\vec{x}^T \vec{x})^{-1}$$

- 之前我们说linear 的bias 大, 是因为不一定原模型就是linear model

- Is there a relationship between X and Y?

- F test: $H_0 : \text{all the } \beta = 0$

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)},$$

- Deciding on importance Variables

- T-test and model selection is different, T-test is想控制type1 error, 也就是在原假设为真但我拒绝了, model selection 是想找出预测最准确的model(focus on testing error), 两个给出的结果不一致, 但在假设正确且数据足够大的情况下, 结果会趋于一致
- 我记得LCM说的是: AIC, BIC都存在E, 但我们只用了一组特定的数据来计算Cp, 所以结果可能不一样

- Model Fit

- 可能会发现, 变量个数越多, R^2 越大, 这是因为我们的 R^2 是用training data 来计算的, 变量越多, 哪怕和谐变量可能并不显著, 也会导致我们的training error 越小, R^2 越大

- Prediction

- model bias(很可能不是linear model)
- reducible error
- irreducible error
- confidence interval

Prediction:

$$\begin{aligned} \text{Var}(\hat{Y}_{\text{new}}) &= \text{Var}(X_{\text{new}}^T \hat{\beta}) = E[X_{\text{new}}^T \hat{\beta} \hat{\beta}^T X_{\text{new}}] \\ &\quad - E[X_{\text{new}}^T \hat{\beta}] E[X_{\text{new}}^T \hat{\beta}]^T \\ &= X_{\text{new}}^T E[\hat{\beta} \hat{\beta}^T] X_{\text{new}} \\ &\quad - X_{\text{new}}^T E[\hat{\beta}] E[\hat{\beta}]^T X_{\text{new}} \\ &= X_{\text{new}}^T (\text{Var}(\hat{\beta})) X_{\text{new}} \end{aligned}$$

*the variance of your estimate
relies on your estimation on β*

- Extensions

- linear model 最重要的两个假设: 1. additive 2. linear
- 什么是additive: X_j 对Y的影响相互独立

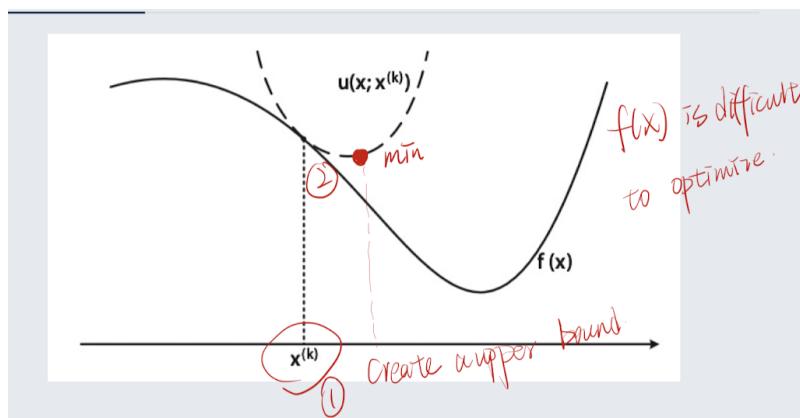
additive
 model $Y = \beta_0 + f(x_1) + \dots + f(x_n)$

- 如果我们remove这个assumption, $\beta_3 \neq 0$ 说明, X_1 对Y的影响还取决于 X_2

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon.$$

- Potential Problems

- non-linearity: residual plot, 不管维度多少, 这个图都非常有用且简单: if residual is not zero --> linear model is not appropriate and try to add some nonlinear parts
- correlation of error terms: ϵ 不独立, 那么 $Var(\hat{\beta}) = \sigma^2(X^T X)^{-1}$, 这个表达式会低估 Var, 所以 confidence interval 会偏窄, 也就是说, 真正的这个区间的把握应该小于 95% --> GEE 可能可以解决这个问题
- non-constant variance of error terms: 我们之前假设 $Var(\sigma^2)$ 是不变的, 但如果他变化呢? 可以变形 $Y \rightarrow \log(Y)/\sqrt{Y}$ (concave function) Y 会变小, 异方差的影响也就会更小一点
 - 有时一些 X 的方差我们是已知的, 这时就需要修正我们的方差,
- outliers: 1. largely deviate from other data 2. 数量很少
 - how to solve:
 - method 1
 - $\omega_1 = 1 \rightarrow$ find outliers
 - change $\omega_i = \frac{1}{|e_i|}$
 - method 2
 - 首先我们用 L1 代替 L2 --> 因为 L2 对数据的变化太敏感
 - 但 L1 不好的地方在于它是 non-differentiable 的
 - 所以我们找一个 L1 的 upper bound 且需要 differentiable



- x_k 已知, 找原曲线上该点的upper bound function, 找这个function的最小区间, --> x_{k+1}
- but how to construct u function?

- formula

$$u(t, t^k) = \frac{1}{2|t^k|} \left(t^2 + (t^k)^2 \right). \quad \text{HT}$$

$$\min_{\beta} u(\beta, \beta^k) = \sum_{i=1}^N \frac{1}{2|[X\beta^k - y]_i|} \left([X\beta^k - y]_i^2 + ([X\beta^k - y]_i)^2 \right).$$

- algorithm

- Set $k = 0$ and initialize with a feasible point β^0
 - **repeat**
 - $w_i^k \leftarrow \sqrt{\frac{1}{2|[X\beta^k - y]|_i}}$
 - Update $\beta^{k+1} \leftarrow \operatorname{argmin}_{\beta} \|(X\beta - y) \odot w^k\|$
 - $k \leftarrow k + 1$
 - **until** convergence
 - **return** β^k

- 注意：不能用residual来判断是否有outlier，因为我们的residual plot和我们用outlier预测出来的model是有关系的，有可能他针对元数据就是outlier，但拟合的时候我们不认为它是outlier
 - high leverage points $\frac{\partial \hat{y}_i}{\partial y_i}$ 非常大，也就是这个点对 y_i 的变化非常敏感
 - model flexibility =sensitivity, 即我对于这个点，仅仅换一个training set的y值，我拟合的这个model变化大不大？
 - 猜想--> $\sum \frac{\partial \hat{y}_i}{\partial y_i}$, how to prove? (if it can cover the special case, then we think it is correct.
eg. linear model's flexibility is the number of parameters)

$$\hat{\beta} = (X^T X)^{-1} X^T y, \quad \hat{y} = X \hat{\beta} = \underbrace{X(X^T X)^{-1} X^T y}_{H},$$

$$\hat{y}_i = h_{ii} y_i + h_{i2} y_2 + \dots + h_{in} y_n, \quad \frac{\partial \hat{y}_i}{\partial y_j} = h_{ij}.$$

$$\sum_{i=1}^{n_{\text{train}}} \frac{\partial \hat{y}_i}{\partial y_j} = \text{tr}(H) = \text{tr}(X(X^T X)^{-1} X^T) = \text{tr}(X^T X (X^T X)^{-1}) = p^n.$$

= number of parameters.

H is symmetric.

(1) H is semi-positive.

$$(2) h_{ii} = h_{11}^2 + h_{22}^2 + \dots + h_{nn}^2 = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2 \geq h_{ii}^2.$$

$$\Rightarrow -1 \leq h_{ii} \leq 1.$$

- flexibility is upper bounded by number of data --> so we need many many data
- but \hat{y}_i is random variable, --> 以上公式random variable, 所以为了获得一个更准确的公式, we need expectation

- example

$$\begin{aligned} & \min_{\beta_0, \beta} \sum_i (y_i - (\beta_0 + \beta x_i))^2 \\ & \text{s.t. } \beta_0 = 1. \end{aligned}$$

如果用之前的 def $\rightarrow \infty$

新 def. $E[\cdot] \rightarrow 1$ correct

MTN $\sum_i (y_i - (\beta_0 + \beta x_i))^2$ s.t. $\beta_0^2 + \beta^2 \leq T$

if $T=0$, $\beta_0 = \beta = 0$, $y_i = 0$

T very large \rightarrow entire space

flexible structure of model constraints

- Collinearity: eg $X_1 = \alpha X_2$
 - 两两之间: pairs
 - 有的不能直观看出<-> multi collinearity<-> VIF: 越大说明越有collinearity

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R^2_{X_j|X_{-j}}},$$

where $R^2_{X_j|X_{-j}}$ is the R^2 from a regression of X_j onto all the other predictors. If $R^2_{X_j|X_{-j}}$ is close to one, then collinearity is present, and so the VIF will be large. *how to detect? $X_j \sim X_{-j}$*

- 但可能没那么好，因为一个一个做我们需要做p次， $\lambda > 0$ ，这应该是个修正，本来如果共线性， $(X^T X)^{-1}$ 是不存在的，在这加了一个 λ 让它存在

$$(X^T X + \lambda I)^{-1} X^T y$$

- VS K-Nearest Neighbors (KNN)
 - non-parametric
- 不能简单认为原数据是non-linear的，就该用non-linear拟合
- Classification
 - supervised learning
 - why not regression
 - qualitative response 不能用量化数字来表示
 - 类别有时是有顺序的，例如severe/moderate/mild，但如果用2, 1, 0来表示，数字说明他们之间的差别是相同的，但实际上不是这样的，我们不知道程度之间的差别到底相不相等，这是无法量化的
 - 我们没办法确保 $\hat{y} \in [0, 1]$
 - LDA(linear discriminant analysis)
 - sigmoid function/softmax active function

$$\text{basic assumptions } p(y_i=1|x_i) = \frac{1}{1 + \exp(-x_i^T \beta)} \in [0, 1].$$

$$p(y_i=0|x_i) = \frac{\exp(-x_i^T \beta)}{1 + \exp(-x_i^T \beta)} = \frac{1}{1 + \exp(x_i^T \beta)}$$

- assumptions: sigmoid function $+ y_i | x_i$ is independent
- why this function? can we choose another function?
 - Yes, for example $\int_0^x N(z) dz$, z follows normal distribution-->probit model
 - The performance are nearly the same
- logistic function: 1. nearly unbiased 2. easy interpretation
- in fact, it is the generalization of linear model

$$\underbrace{\log \frac{P(y_i=1|x_i)}{P(y_i=0|x_i)}}_{\text{log odds ratio.}} = \frac{\log}{\log \exp(x_i^T \beta)} = x_i^T \beta.$$

- we have two perspectives to understand

- method 1

maximum likelihood function

$$\prod_{i=1}^n P(y_i|x_i, \beta) = \prod_{i=1}^n P(y_i=1|x_i, \beta) \cdot \prod_{i=1}^n P(y_i=0|x_i, \beta)$$

$$\max(\log P(D|\beta)) = \max \left(\sum_{i=1}^n \log P(y_i|x_i, \beta) \right)$$

$$= \max \left(\sum_{i=1}^n y_i x_i^T \beta - \log \left(\prod_{i=1}^n \exp(x_i^T \beta) \right) \right)$$

unified formula.

$$P = \frac{\exp(x_i^T \beta)}{\prod_{i=1}^n \exp(x_i^T \beta)} = \frac{1}{\prod_{i=1}^n \exp(x_i^T \beta)} P(y_i=1|x_i, \beta)$$

- method 2

$$P_i = P(y_i=1|x_i)$$

$$1-P_i = P(y_i=0|x_i)$$

$$P = P_i^{y_i} (1-P_i)^{1-y_i}$$

why correct?
when $y_i=1 \rightarrow P_i$
when $y_i=0 \rightarrow 1-P_i$

And this is the ~~MLE~~ loss function.

$$\Rightarrow MLE = \prod_{i=1}^n P_i^{y_i} (1-P_i)^{1-y_i} = L(\theta) = P(D|\beta)$$

$$\max(\log L(\theta)) = \max \left(\log \prod_{i=1}^n P_i^{y_i} (1-P_i)^{1-y_i} \right)$$

$$= \max \left(\sum_{i=1}^n y_i \log P_i + (1-y_i) \log (1-P_i) \right).$$

cross entropy.
 $-\log P(D|\beta)$ cross entropy loss function.

- why equal?

why equal?

$$\text{if } y_i=1, \text{ method 2} = \log P_i = \log \frac{\exp(x_i^T \beta)}{1+\exp(x_i^T \beta)}$$

$$= x_i^T \beta - \log(1+\exp(x_i^T \beta)) = \text{method 1.}$$

$$\text{if } y_i=0, \text{ method 2} = \log(1-P_i) = \log \frac{1-\exp(-x_i^T \beta)}{1+\exp(-x_i^T \beta)}$$

$$= \log \frac{1}{1+\exp(x_i^T \beta)}$$

$$= -\log(1+\exp(x_i^T \beta)) = \text{method 1.}$$

- why we need the method 2? Because it is easier to generalize

$$P_1 y_1 P_2 y_2 P_3 y_3 \dots P_1 + P_2 + P_3 = 1. \quad y_1 + y_2 + y_3 = 1 \quad y_i = 0/1 \\ \Rightarrow \text{only one of } y_i = 1, \text{ else } y_i = 0.$$

- how to solve the MLE?

- 一般来说我们求导计算

$$L(\beta) = \sum_i y_i x_i^\top \beta - \log(1 + \exp(x_i^\top \beta))$$

$$\frac{\partial L}{\partial \beta} = \sum_i y_i x_i - \left(\frac{\exp(x_i^\top \beta)}{1 + \exp(x_i^\top \beta)} \right) x_i \cdot P_i = \frac{1}{1 + \exp(-x_i^\top \beta)}$$

$$= x_i(y - p)$$

$$\frac{\partial^2 L}{\partial \beta^2} = \sum_i \frac{x_i}{(1 + \exp(-x_i^\top \beta))^2} \cdot \exp(-x_i^\top \beta) \cdot (-x_i^\top) \quad \text{not compatible.}$$

how to solve? iteration

- 但这里我们不能解决二阶导, how to solve?

Procedure:

$$\frac{\partial^2 L}{\partial \beta^2} = \sum_i - \underbrace{P_i(1-P_i)}_{\text{diag}(W)} x_i x_i^\top.$$

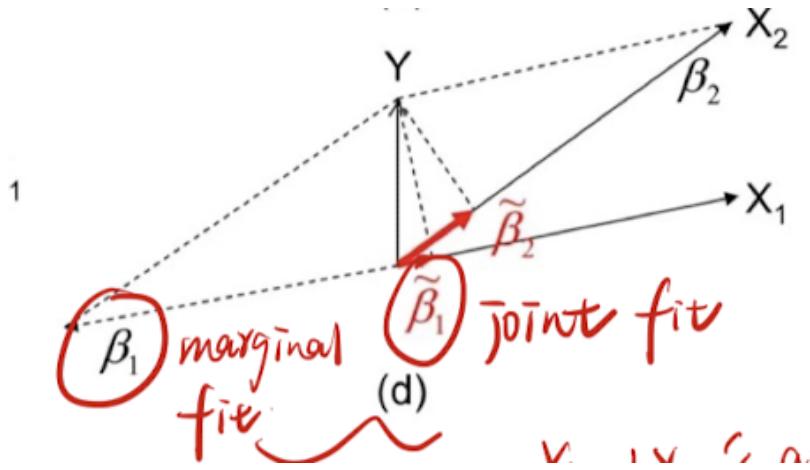
We let $W = \text{diag}(P_1(1-P_1), P_2(1-P_2), \dots, P_n(1-P_n))$

$$\beta^{\text{new}} = \beta^{\text{old}} - H^{-1} \nabla L |_{\beta^{\text{old}}}$$

why correct? $\min_x f(x) \quad g(x) = \nabla f(x) = 0$, but I can't solve it.
 \Rightarrow 1st approximation of $g(x)$: $g(x) \approx g(x_0) + g'(x_0)(x - x_0) = 0$

$$x = x_0 - \underbrace{g'(x_0)^{-1}}_{\text{H}} \underbrace{g(x_0)}_{\nabla L|_{\beta^{\text{old}}}}$$

- 我们得到的结果会依赖于 β_0 吗
 - 不会, 因为这是一个convex problem-- 二阶导大于等于0
- example and paradox
 - 我们得到的系数都是在有其他系数的影响下得到的系数, 也就是影响程度, 即joint fit
 - 如果仅有一个X因子, 才是不condition on other factors--> marginal fit
 - depending on the angle of these two factors
 - marginal fit 是平行四边形分解, 而joint fit 是有Y向两个X作垂线



- 两个系数都正确，取决于我们如何解释这个系数，joint fit即是conditional的
 - Multi-class
 - 我们首先分析一下两个名词：
 - generative model：什么是生成式模型，就是我在给定这个模型的结构和参数之后，可以直接生成数据，最简单的就是：如果我知道了一个函数f的分布，那么我可以不需要任何别的参数就可以生成符合这个分布函数的数据-->从这个角度来看，logistic Reg不是生成式模型，而是判别式模型，discriminative model(给定了模型的结构和参数，还必须知道X才能生成符合条件的Y，如果没有X，就什么也做不了)
 - 什么模型是生成式模型呢？我们想判断Y到底是属于哪一类，在已知这两个分布之后，对K类中的每一类我们都能算出一个概率，那个概率最大，我们就认为属于哪一类
- 如果已知 $x|y=k \sim N(\mu_k, \Sigma)$ $P(y=k)$
 我们能得到 $y_i \leftarrow P(y=k)$
 $x_i \leftarrow P(x=m) = \sum_k P(y=k) \cdot P(x=m|y=k)$
 $\underbrace{P(y_k|x)}_{\downarrow} = \frac{P(x,y_k)}{\sum_k P(x,y_k)} = \frac{P(x|y_k)P(y_k)}{\sum_k P(x|y_k)P(y_k)}$
 Then we can get (x_i, y_i)
- 这与判别法不同的在于，判别式模型有一条明显的分界线，在分界线的哪边就属于哪一类，等于是在整个数据集上训练了一个函数，而生成式模型是对每一类都训练了一个函数，即每一类都有一个概率分布
 - LDA theory
 - theory: $\delta_k(x)$ is called discrimination function

Assume: k classes. $X|Y=k \sim N(\mu_k, \Sigma)$ have the same covariance matrix

$$\begin{aligned} \Rightarrow p_k(x) &= P(X|Y=k) \\ &= \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) \right) \\ \Rightarrow p_k(x) &= P(Y=k|X) = \frac{P(X|Y=k)P(Y=k)}{\sum_i P(X|Y=i)P(Y=i) \pi_i} \\ &= \frac{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}} \exp(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k)) \pi_k}{\sum_i (2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}} \exp(-\frac{1}{2}(x - \mu_i)^T \Sigma^{-1} (x - \mu_i)) \pi_i} \\ &= \frac{\pi_k \exp(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k))}{\sum_{i=1}^k \pi_i \exp(-\frac{1}{2}(x - \mu_i)^T \Sigma^{-1} (x - \mu_i))} \end{aligned}$$

我们的目标是找出 $p_k(x)$ 最大的 k , 即 $\log p_k(x)$ 最大, 分母都相同

因此只看分子的大小, 即 $p_k(x) > p_i(x)$

$$\begin{aligned} \log p_k(x) - \log p_i(x) &= \log \pi_k - \frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) \\ &= \log \pi_k - \frac{1}{2} x^T \Sigma^{-1} x - \cancel{x^T \Sigma^{-1} \mu_k} - \cancel{\mu_k^T \Sigma^{-1} x} + \cancel{\mu_k^T \Sigma^{-1} \mu_k} \\ &= \log \pi_k - \frac{1}{2} x^T \Sigma^{-1} x + \cancel{-\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k} \\ \log p_k(x) - \log p_i(x) &= \log \pi_k - \cancel{\frac{1}{2} x^T \Sigma^{-1} x + \mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k} \\ &\quad - \log \pi_i - \cancel{\frac{1}{2} x^T \Sigma^{-1} x - \mu_k^T \Sigma^{-1} x + \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k} \\ \text{So } S_k(x) &= \mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k \\ &= (\mu_k^T \Sigma^{-1} \mu_k)^{\frac{1}{2}} x + \beta_0 = \beta_1^T x + \beta_0 - \log p_k(x) \end{aligned}$$

logistic function!

If $S_k(x) = S_i(x) \Rightarrow$ we can solve the boundary.

- but how to compute the parameters eg. π_k, μ_k
 - μ_k 就是这一类的均值

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i;$$

- σ 怎么计算, 首先我们认为所有的都是一个方差, 所以必须用到所有的数据, 但每一组的均值不一样, 所以还是要分组求, 注意除得是什么, 是自由度, 因为我们有 k 个均值, 因此我们的自由度是 $n - k$ (所以之前的方差为什么用 $n - 1$ 来除也是一样的道理)

$$\hat{\Sigma} = \frac{1}{n - k} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T,$$

variance of special case k .
df: the number of k classes

- π_k 就是 Y 的概率分布, 直接用经验分布就好了
- 还有一个 LDA(latent dirichlet allocation): topic model: 即给一个 document, 他可以分类这篇 document 有 40% science, 50% nature and 10% finance 类似这种
- Comparison of logistic Reg and LDA
- same:

Same: $\log \frac{P(Y=1|x)}{P(Y=2|x)} = \beta_1^T x + \beta_0$ all linear func.

- difference:
 - 1. LDA assume the structure so we can just calculate, but logistic Reg don't have the closed form, so we use Newton + MM
 - 2.还记得吗, 对于logistic Regression, 因为我们用了迭代, 当true optimal value is $+\infty$ 时, 我们的迭代是不能收敛的, 也就是说当两个数据集时完全分开的时候, β 必须趋于无穷, 概率才能为1, 所以是存在numeric problem 的, 但LDA没有这种问题
 - 3.LDA is a special case of logistic Reg

- how to evaluate the method?

- 很自然的想法是misclassification error, 也被称为averaged misclassification error, or zero-one loss function, 但这不是个好的方法, 比如下面的例子, 如果我们完全不分类, 仅仅把所有的都认为是no, 我们的正确率会更好, 另一种解释是我们更多时候更关注class one, 所以我们要换个别的方法来评估

Predicted default status

True		No	Yes	Total	
Predicted	No	9644	23	9667	Total
	Yes	252	81	333	
Total	9896	104	10000		cla

- 首先我们定义: FP,FN,TP,TN,第一个字母代表我们预测的是否正确, 第二个字母代表我们预测的是什么

$$\begin{aligned}
 \text{FP} &= p(\hat{y}_i=1 | y_i=0) \\
 \text{FN} &= p(\hat{y}_i=0 | y_i=1) \\
 \text{TP} &= p(\hat{y}_i=1 | y_i=1) \quad \text{sensitivity: ability to identify positive classes} \\
 \text{TN} &= p(\hat{y}_i=0 | y_i=0) \quad \text{specificity}
 \end{aligned}$$

- 在每个阈值下, 我们都会存在一个sensitivity 和specificity, 这个阈值是比如大于0.5还是大于0.2我们把它归类为1, 随着阈值变小, TP和FP都应该升高, 因为归类为positive的变多, 也就是说, 对于一个model, 我们有这样的一个随着阈值变化的图, 我们称为ROC curve, 这条曲线下面的面积成为AUC, 我们把这个数值认为是衡量该模型表现的一个指标
- 我们还有一种比较model好坏的, 去比较loss function, 两者不同是我们这个AUC是没有假设任何模型的, 但logistic是假设了一些分布, 比如伯努利
- QDE

- 我们上面的LDA是假设了是线性，也就是大家的方差是一样的，但还有另外一种，我嫩不认为他是一样的，但这样会有很多项就不能消掉了，

$$\delta_k(x) = -\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) - \frac{1}{2} \log(|\Sigma_k|) + \log \pi_k$$

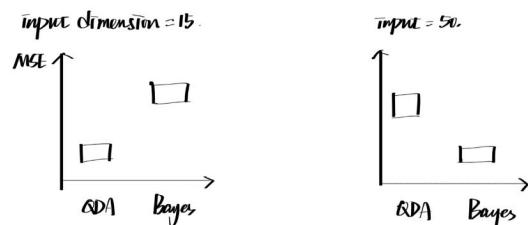
$$= -\frac{1}{2} x^T \Sigma_k^{-1} x + x^T \Sigma_k^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma_k^{-1} \mu_k - \frac{1}{2} \log(|\Sigma_k|) + \log \pi_k$$

This is a quadratic function of x .

- 注意：如果variances相近，LDA更好，不想尽，QDA更好？！！！！！NO，这个说法是错误的

- 好不好取决于MSE，也就是bias和variance，并不取决于assumptions

QDA : $X|Y \sim N(\mu_Y, \Sigma_F)$.
 $D_{\text{train}}(500)$, $D_{\text{test}}(50)$: $X|Y \sim N(\mu_Y, \Sigma_F)$; Σ_F is diagonal $\Rightarrow X|Y$ is independent
 Naive Bayes



why? If dimension = 50. QDA. we must estimate 500 parameters

but I only have 500 data points
 \Rightarrow MSE 很大

But if Naive Bayes. \Rightarrow So we just only estimate $p=50$ parameters

认为 X 之间 independent \Rightarrow MSE 小

So, 一个 model 不仅并不取决于 model assumption 是否正确。

而是取决于 Bias & Variance trade-off

- 但是QDA需要更多的数据，如果数据太少，那个LDA可能更合适