

Each observation is characterized by different variability

Example: each obs. is an average of different number of points

Recall: for iid Y_i , $\text{Var}\left(\frac{1}{N} \sum_{i=1}^N Y_i\right) = \frac{1}{N} \text{Var}(Y_i)$

In linear regression: $Y = X\delta + \epsilon$, $\text{Var}(\epsilon) = \sigma^2 \Sigma$, $\Sigma = \begin{pmatrix} \frac{1}{n_1} & \frac{1}{n_2} & \dots & \frac{1}{n_N} \end{pmatrix}$

Change of notation: $Y = X\delta + \epsilon$, $\text{Var}(\epsilon) = \sigma^2 W^{-1}$, $W = \begin{pmatrix} n_1 & n_2 & \dots & n_N \end{pmatrix}$

Linear transformation: $W^{1/2}Y = W^{1/2}X\delta + W^{1/2}\epsilon$

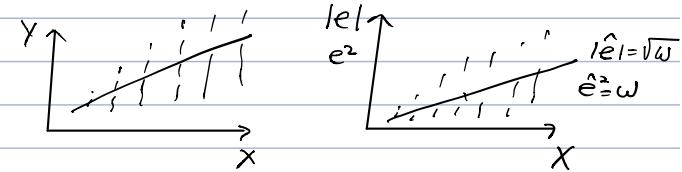
$$\text{Var}(\epsilon_w) = \text{Var}(W^{1/2}\epsilon) = \sigma^2 W^{1/2} \Sigma W^{1/2} = \sigma^2 \Sigma W^{1/2} W^{1/2} = \sigma^2 I$$

Change in loss function

Minimizing squared error on the transformed scale:

$$\min_{\delta} (Y_w - X_w \delta)^T (Y_w - X_w \delta) = (W^{1/2}Y - W^{1/2}X\delta)^T (W^{1/2}Y - W^{1/2}X\delta) = (Y - X\delta)^T W (Y - X\delta) = \sum_{i=1}^N w_i (y_i - x_i \delta)^2$$

$\Rightarrow \hat{\delta} = (X_w^T X_w)^{-1} X_w^T Y_w = (X^T X)^{-1} X^T W Y - \text{averages with larger } n_i \text{ have more weight}$



Example $\text{Var}(y_i | x_i) = f(x_i)$

Approach: (1) Fit OLS regression

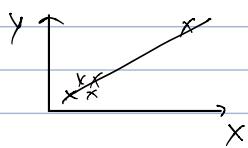
(2) Find a relationship between $|e_i|$ or e_i^2 and x_i

(3) Use the relationship to quantify weights

(4) Refit regression with the weights

(5) Optional: iterate (1)-(4)

Observations with high leverage



Recall $\hat{Y} = \underbrace{X(X^T X)^{-1} X^T}_H Y$

H : hat matrix

Also outliers -
see hw 1

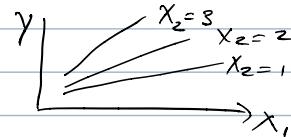
Leverage statistic for obs. i : i th diagonal element of H .

2-variable regression: $h_{ii} = \frac{1}{N} + \frac{(x_{ii} - \bar{x})^2}{\sum_{j=1}^N (x_{ij} - \bar{x})^2}$

Expressing non-linearity

Polynomial regression: $y = \delta_0 + \delta_1 x + \delta_2 x^2 + \dots + \delta_{p-1} x^{p-1} + \epsilon$

Regression with higher order terms: $y = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \delta_3 x_1 x_2 + \epsilon$



-surface with curvature

Price for more flexible models: more parameters

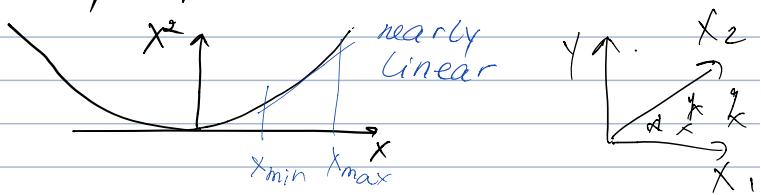
Multicollinearity

$\hat{\delta} = (X^T X)^{-1} X^T Y$ cannot invert $X^T X$ when $\text{rank}(X) < p$

- X has linearly dependent columns
- $N < P$

Example: higher order terms

Solution: center and scale



Bias-variance tradeoff

Until now assumed that the linear model is true. How to select the right model? what is the price for the wrong model?

Consider loss functions again:

$J(\delta) = \sum_{i=1}^N (y_i - \hat{h}_\delta(x_i))^2$ - squared loss/squared error: evaluate $h_\delta(x)$ on current data

Generalization error: performance on future, independent test data
= prediction error

Error_{test} = $E[J(\delta)|\{x_i, y_i\}]$ - expected loss, ave. over all future data, given training set $\{x_i, y_i\}$ - more meaningful

Error = $E[J(\hat{\delta})]$ - expected loss, ave. over all future data and all training sets
- easier mathematically

Prediction for $x = x_0$ squared loss

$$\begin{aligned} \text{Error}(x_0) &= E[Y - \hat{h}_\delta(x_0)]^2 = E[Y - h(x_0) + h(x_0) - E[\hat{h}_\delta(x_0)] + E[\hat{h}_\delta(x_0)] - \hat{h}_\delta(x_0)]^2 \\ &= \underbrace{E[Y - h(x_0)]^2}_{\text{irreducible error}} + \underbrace{E[h(x_0) - E[\hat{h}_\delta(x_0)]]^2}_{\text{bias}^2} + \underbrace{E[E[\hat{h}_\delta(x_0)] - \hat{h}_\delta(x_0)]^2}_{\text{variance}^2} \\ &\quad \text{- cross-products cancel out} \qquad \text{model mis-specification} \qquad \text{variance of the sampling distribution of parameter estimates} \end{aligned}$$

$$[\text{ave model bias}]^2 + [\text{ave estimation bias}]^2$$

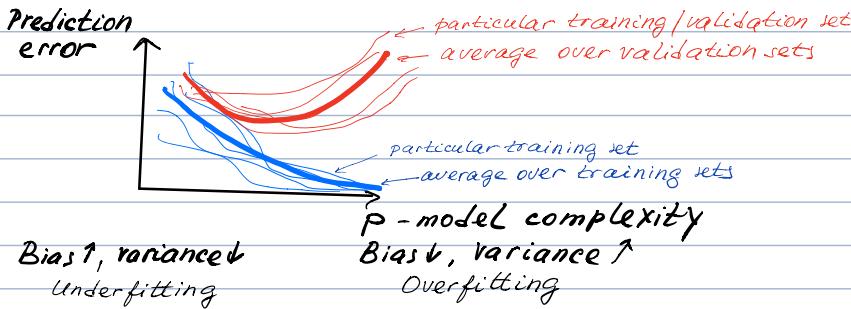
Specific to linear regression:

$$\begin{aligned} \text{Error}(x_0) &= \sigma^2 + E[h(x_0) - E[\hat{h}_\delta(x_0)]]^2 + \underbrace{\|X(X^T X)^{-1} x_0\|_2^2 \sigma^2}_{\rightarrow E[(x_0^T (X^T X)^{-1} X^T Y) - x_0^T (X^T X)^{-1} X^T Y]^2} \\ &= E\{(x_0^T (X^T X)^{-1} X^T (E[Y] - Y))\}^2 \\ &= \|X(X^T X)^{-1} X^T \text{Var}(Y) [X^T (X^T X)^{-1} X^T]\|^2 \\ &= \|X(X^T X)^{-1} X^T\|_2 \sigma^2 \end{aligned}$$

How to estimate the error in practice? Average over future data

$$\frac{1}{N} \sum_{i=1}^N \text{Error}(x_i) = \hat{\sigma}_\epsilon^2 + \frac{1}{N} \sum_{i=1}^N [h(x_i) - E[h_\theta(x_i)]]^2 + \frac{P}{N} \hat{\sigma}_\epsilon^2$$

Training data	$P \uparrow \Rightarrow \downarrow$	\downarrow	\downarrow	- apparent (in-sample) estimates
Future data	$P \uparrow \Rightarrow -$	\downarrow	\uparrow	- true (independent) estimates



Goal: minimize prediction error on the validation set

In Linear regression, model selection / model building
= variable selection

Steps of variable selection $N \gg p$

- | |
|--|
| |
| |
| |
| |
- | |
|--|
| |
| |
| |
| |
- model fitting / training
 - model development
 - model evaluation
- X X

(1) Explore the space of candidate models

(2) Evaluate the candidate models on the dev. set, keep the best performing

(3) Evaluate the selected model on the evaluation set

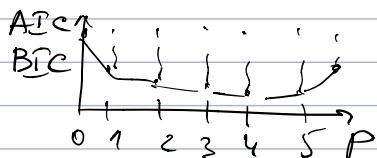
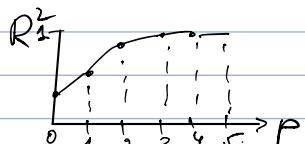
Exploring the model space in (1):

- Exhaustive search of $\binom{p}{k}$ predictors
 - Forward stepwise selection
 - Backward stepwise selection
- many heuristics exist

Example: forward stepwise selection

* Denote M_0 the null model with no predictors

- For $k = 0, \dots, p-1$:
 - consider all $p-k$ models that augment M_k with one predictor
 - $M_k \leftarrow$ the best among $p-k$ models
- R^2
AIC
BIC
- Select the model among M_0, \dots, M_p that performs best on dev set



$$AIC = \frac{SSE}{MSE_P} + 2P$$

$$BIC = \frac{SSE}{MSE_P} + \log N P$$

For Normal distribution \Rightarrow AIC and BIC

$$-2 \log L = \frac{\sum_{i=1}^N (x_i - \hat{x})^2}{\sigma^2}$$

σ^2 estimated by variance from the saturated model

are penalized log-L criteria

Data-poor situation: cross-validation $N \gg p$

Iterate partitioning into train and dev sets

dev \rightarrow



Problem: each iteration does not always have the same best model.

"Consensus" model may not perform as well

Frequently used as an intermediate parameter tuning step