This is called model selection by oos error.

This is one such method. Other methods are "analytical"
i.e. they rely on statistical/prob. models — we won't talk about
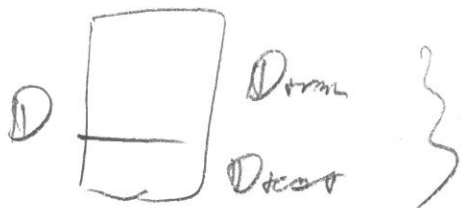those!    AIC, BIC, Cp, etc...

probably won't get to
it this semester
unfortunately.

[ DEMO ]

$\not{\cancel{}}$ midterm 6
$\not{\cancel{}}$ FINAL

There is an extra step in validation that is usually employed.
Remember before vacation we were looking at one model



$D \begin{array}{|c|} \hline D_{train} \\ \hline D_{test} \\ \hline \end{array}$    $\Bigg\}$    random 80% / 20% split of $n$ observations

what if you get an
"unlucky" split?

e.g. all the "weird" observations in $D_{test}$
$\Rightarrow$ your $\underset{\wedge}{\text{error}}$ estimate is much higher
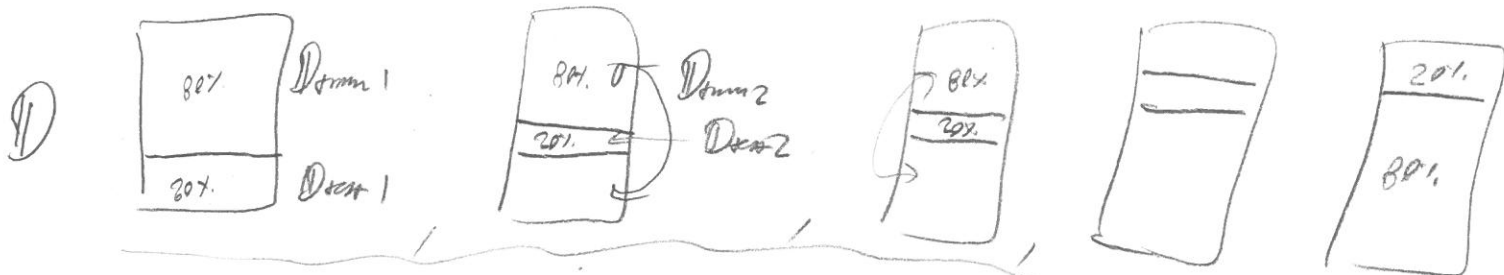than it should be

all the "weird" observations in $D_{train}$
$\Rightarrow$ your error is much lower
" " " " "

$\Rightarrow$ Variance
in the estimate
of future
performance

Solution ?    "K-fold Cross Validation" (CV)

① 

Now, each observation is in the $D_{test}$ once.      $K=5$

Protocol:

① fit $g_k = A(\mathcal{H}, D_{train,k})$

② Save $\vec{\hat{y}}_k = g_k(X_{test,k})$

③ Repeat 1-2 for each fold

④ Concatenate vertically $\vec{\hat{y}}_{CV} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_K \end{bmatrix}$

full $\vec{y}$ since each obs. represented across the K folds

⑤ Compute $oose = error(\vec{y}, \vec{\hat{y}}_{CV})$

↑ Not dependent on a split

Still: maybe those 5 splits were idiosyncratic

⇒ repeat K-fold CV, do the whole procedure many times and avg.

E) CV gives you a lower variance estimate of future performance.

How to choose K? Let's see what happens with different values of K

No Theory exists to pick opt. value ⇒ Open problem in Stat!!

K=2    50% train, 50% test

Since I have 50% test ... I get a very nice estimate of error but on a model which has more est. error! than the final model
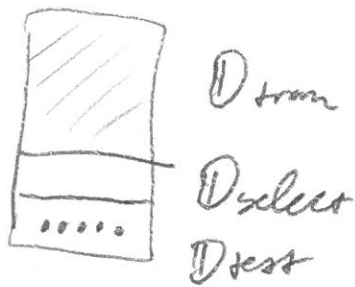
K=5    80% train, 20% test

K=n    leave one out CV" (LOOCV)

Totally random est. of error but a model with low est. error

Optimize tradeoff est. error of g est. error of oose

# CV for model selection

Bootstrap Validation
sample 4 rows from $D$ "oob"
with replacement ... validate on "oob"
"bag"

## For HW



$D_{train}$

$D_{select}$

$D_{test}$

There are two CV loops now!

e.g. $K = 5$

Inner K $(k_i)$

Outer K $(k_o)$

      1      2      3      4      5



## Protocol

① for each outer loop: $(k_o)$

② for each inner model $j \in \{1, ..., M\}$

③ for each inner loop $(k_i)$

① $\forall \{k_i, k_o\}$: $F_j(Z_j, D_{train, k_i, k_o})$

② compute $\vec{\hat{y}}_{j, i, k_i, k_o} = \hat{g}_{j, i, k_i, k_o}(D_{select, k_i, k_o})$

③ Repeat 1-2 for all models $j \in \{1, ..., M\}$

④ Repeat 1-2 for all inner folds $k_i \in \{..., 5\}$

⑤ Concatenate

$$\vec{\hat{y}}_{j, k_o} = \begin{bmatrix} \vec{\hat{y}}_{j, 1, k_o} \\ \vdots \\ \vec{\hat{y}}_{j, 5, k_o} \end{bmatrix}$$

⑥ Select best model

$$j^*_{k_o} = \operatorname{argmin}\{\text{oose}_{j, k_o}, ..., \text{oose}_{M, k_o}\}$$

⑦ Repeat 1-6 for $k_o \in \{... 5\}$

⑧ set $\vec{\hat{y}} = \begin{bmatrix} \hat{y}_{j^*_1, 1} \\ \vdots \\ \hat{y}_{j^*_5, 5} \end{bmatrix}$

⑨ Estimate $\text{oose} = \operatorname{error}(\vec{\hat{y}}, \vec{y})$

⑩ Repeat steps 1-6 to build final model $g$ without $D_{test}$ (only inner CV loop)