# Machine Learning from Data

Lecture 13: Spring 2021

# Today's Lecture

- Validation and Model Selection
  - Validation Set
  - Model Selection
  - Cross validation

# Regularization (Recap)

Regularization combats the effects of noise by putting a leash on the algorithm.
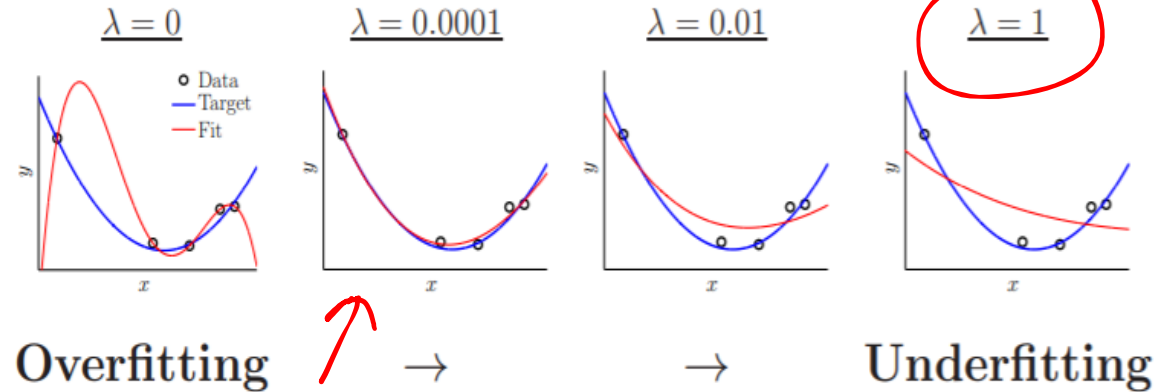
$$E_{\text{aug}}(h) = E_{\text{in}}(h) + \frac{\lambda}{N}\Omega(h)$$

$\Omega(h) \rightarrow$ smooth, simple $h$

noise is rough, complex.

Different regularizers give different results

can choose $\lambda$, the **amount** of regularization.



|  $\lambda = 0$  |  $\lambda = 0.0001$  |  $\lambda = 0.01$  |  $\lambda = 1$  |

**Overfitting** $\rightarrow$ $\rightarrow$ **Underfitting**

Optimal $\lambda$ balances approximation and generalization, bias and variance.

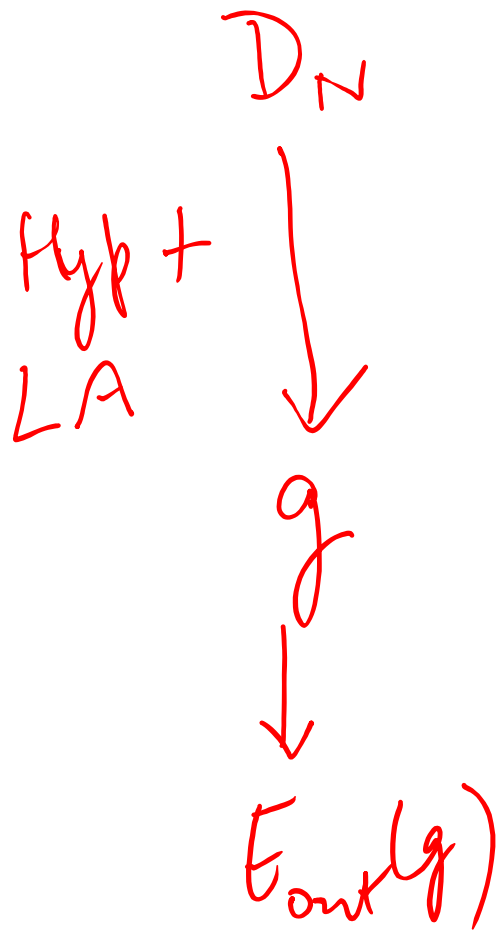$$E_{\text{out}}(g) = E_{\text{in}}(g) + \text{overfit penalty}$$

VC bounds this using a complexity error bar $\Omega(\mathcal{H})$

regularization estimates this through a heuristic complexity penalty $\Omega(g)$

# Validation

Validation goes directly for the jugular:

$$E_{\text{out}}(g) = E_{\text{in}}(g) + \text{overfit penalty}.$$

validation estimates this directly

In-sample estimate of $E_{\text{out}}$ is the Holy Grail of learning from data.

$D_N$

Hyp +
LA

$\downarrow$

$g$

$\downarrow$

$E_{out}(g)$

$\boxed{x}$

$e(g(x), y)$

$e(g)$

$E[e(g)] = E_{out}(g)$

$\underline{Var(e(g))} \uparrow$

Test data set

$x_1 \quad x_2 \cdots \quad x_K$

$\downarrow \qquad \downarrow \qquad\qquad \downarrow$

$e_1 \quad e_2 \qquad e_k$

$E_{test} = \frac{1}{K} \sum_{k=1}^{K} e_k$

$E[E_{test}] = E\left[\frac{1}{K} \sum_{k=1}^{K} e_k\right]$
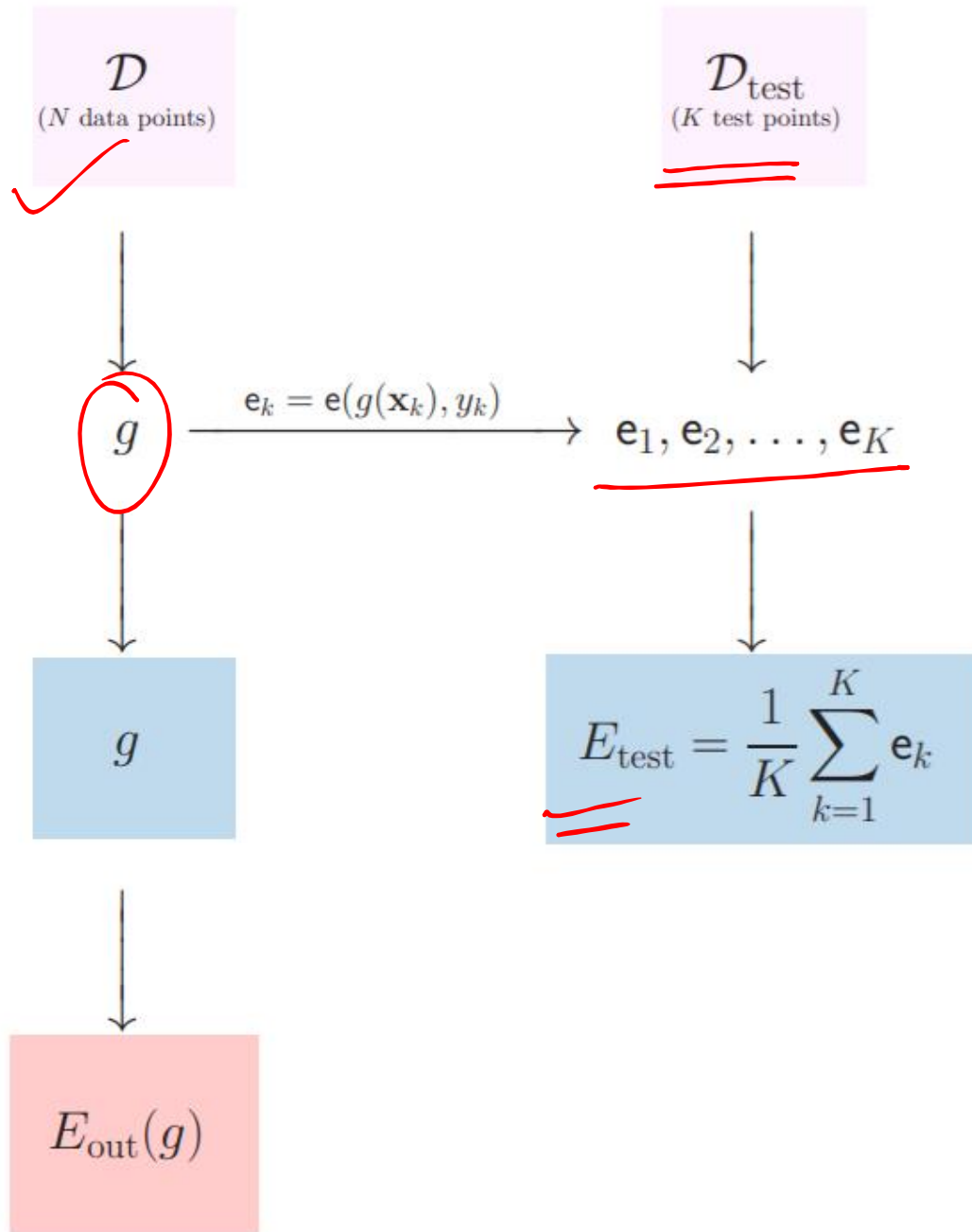
$= \frac{1}{K} \left(\sum_{k=1}^{K} E(e_k)\right)$

$= \frac{1}{K} E_{out}(g)$

$$\text{Var}(E_{test}) = \text{Var}\left[\frac{1}{K}\sum_{k=1}^{K}e_k\right]$$

$$= \frac{1}{K^2}\text{Var}\left[\sum_{k=1}^{K}e_k\right]$$

$$\therefore K \cdot \text{Var}(e_k) = \frac{1}{K}\text{Var}(e_k)$$

$$\text{Var}(E_{test}) = \frac{1}{K^2}$$

$$E_{out} \simeq E_{test} \pm \frac{1}{\sqrt{K}}\sqrt{\text{Var}(e_k)}$$

$\mathcal{D}$

($N$ data points)

$\mathcal{D}_{\text{test}}$

($K$ test points)

$g$ $\quad\xrightarrow{\mathrm{e}_k = \mathrm{e}(g(\mathbf{x}_k), y_k)}\quad$ $\mathrm{e}_1, \mathrm{e}_2, \ldots, \mathrm{e}_K$

$g$

$$E_{\text{test}} = \frac{1}{K}\sum_{k=1}^{K}\mathrm{e}_k$$

$E_{\text{out}}(g)$

$E_{\text{test}}$ is an estimate for $E_{\text{out}}(g)$

$$\mathbb{E}_{\mathcal{D}_{\text{test}}}[\mathrm{e}_k] = E_{\text{out}}(g)$$

$$\mathbb{E}[E_{\text{test}}] = \frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\mathrm{e}_k]$$

$$= \frac{1}{K}\sum_{k=1}^{K}E_{\text{out}}(g) = E_{\text{out}}(g)$$

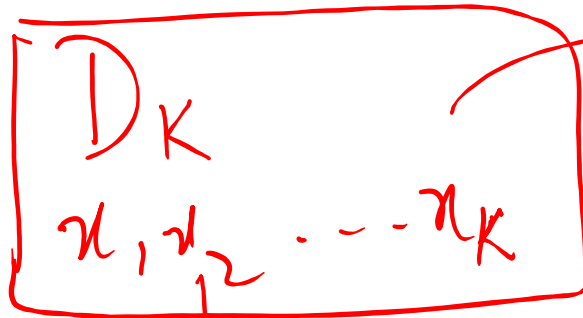$\mathrm{e}_1, \ldots, \mathrm{e}_K$ are *independent*

$$\mathrm{Var}[E_{\text{test}}] = \frac{1}{K^2}\sum_{k=1}^{K}\mathrm{Var}[\mathrm{e}_k]$$

$$= \frac{1}{K}\mathrm{Var}[e]$$

decreases like $\frac{1}{K}$

bigger $K \implies$ more reliable $E_{\text{test}}$.

$D_N$

$D_K$
$x_1, x_2 \cdots x_K$

$\rightarrow$ Validation set

Training set

$D_{N-K}$

$g^-$ (deficient hypothesis)

$E_{out}(g^-)$

$e_1(g^-), e_2(g^-) \cdots e_K(g^-)$

$$E_{val} = \frac{1}{K} \sum_{k=1}^{K} e_k(g^-)$$

$$E[E_{val}] = E_{out}(g^-)$$

$$Var[E_{val}] = \left(\frac{1}{K}\right) Var\left(e(g^-)\right)$$

# The Validation Set



$\mathcal{D}$
($N$ data points)

$\mathcal{D}_{\text{train}}$
($N - K$ training points)

$\mathcal{D}_{\text{val}}$
($K$ validation points)

no data

$g^-$

$\mathbf{e}_k = \mathbf{e}(g^-(\mathbf{x}_k), y_k)$

$\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_K$

$g$

$E_{\text{val}} = \dfrac{1}{K} \sum_{k=1}^{K} \mathbf{e}_k$

$K = N$

$E_{\text{out}}(g^-)$

$E_{\text{val}}$ is an estimate for $E_{\text{out}}(g^-)$

$$\mathbb{E}_{\mathcal{D}_{\text{val}}}[\mathbf{e}_k] \;=\; E_{\text{out}}(g^-)$$

$$\mathbb{E}[E_{\text{test}}] \;=\; \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\mathbf{e}_k]$$

$$= \; \frac{1}{K} \sum_{k=1}^{K} E_{\text{out}}(g^-) = E_{\text{out}}(g^-)$$

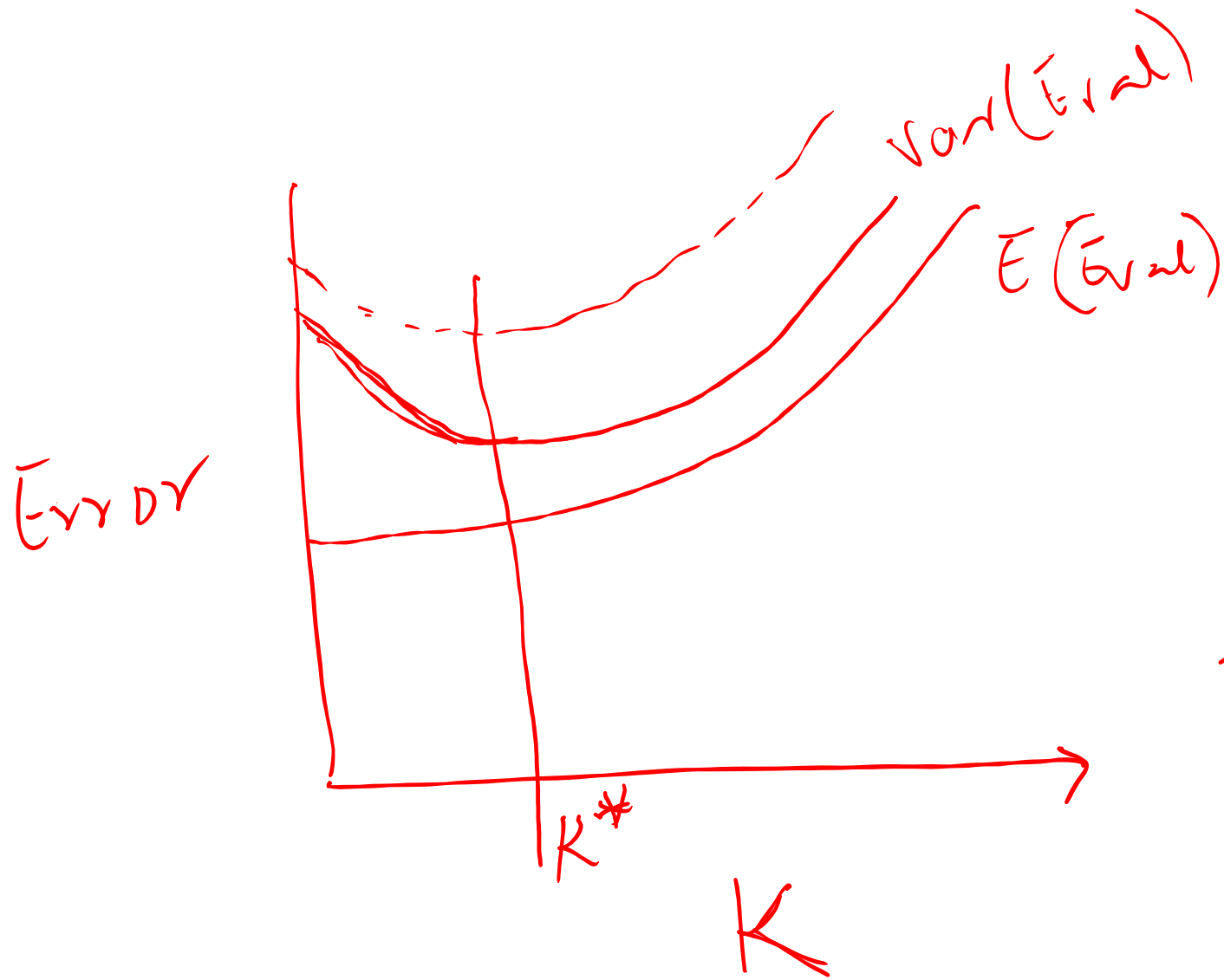$\mathbf{e}_1, \ldots, \mathbf{e}_K$ are *independent*

$$\text{Var}[E_{\text{val}}] \;=\; \frac{1}{K^2} \sum_{k=1}^{K} \text{Var}[\mathbf{e}_k]$$

$$= \; \frac{1}{K} \text{Var}[\mathbf{e}(g^-)]$$

decreases like $\frac{1}{K}$
depends on $g^-$, not $\mathcal{H}$
bigger $K \implies$ more reliable $E_{\text{val}}$?

Var(Eval)

E(Eval)

Error

$K^*$

$K$

$g^- \longrightarrow$ worse

$e(g^-)$

In practice,

$K^* \simeq \dfrac{N}{5}$ or

$20\%$

# Restoring $\mathcal{D}$
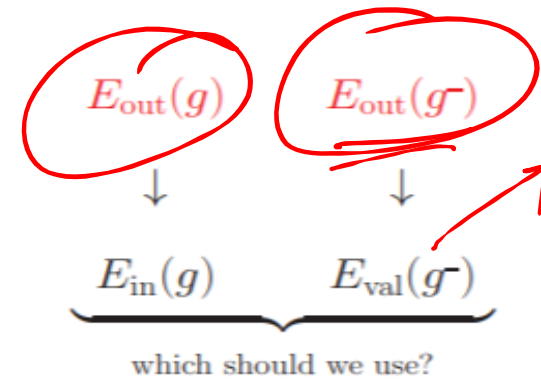


Behind closed doors

**Primary goal:** output best hypothesis.

$g$ was trained on *all* the data. → $N$

**Secondary goal:** estimate $E_{\text{out}}(g)$.
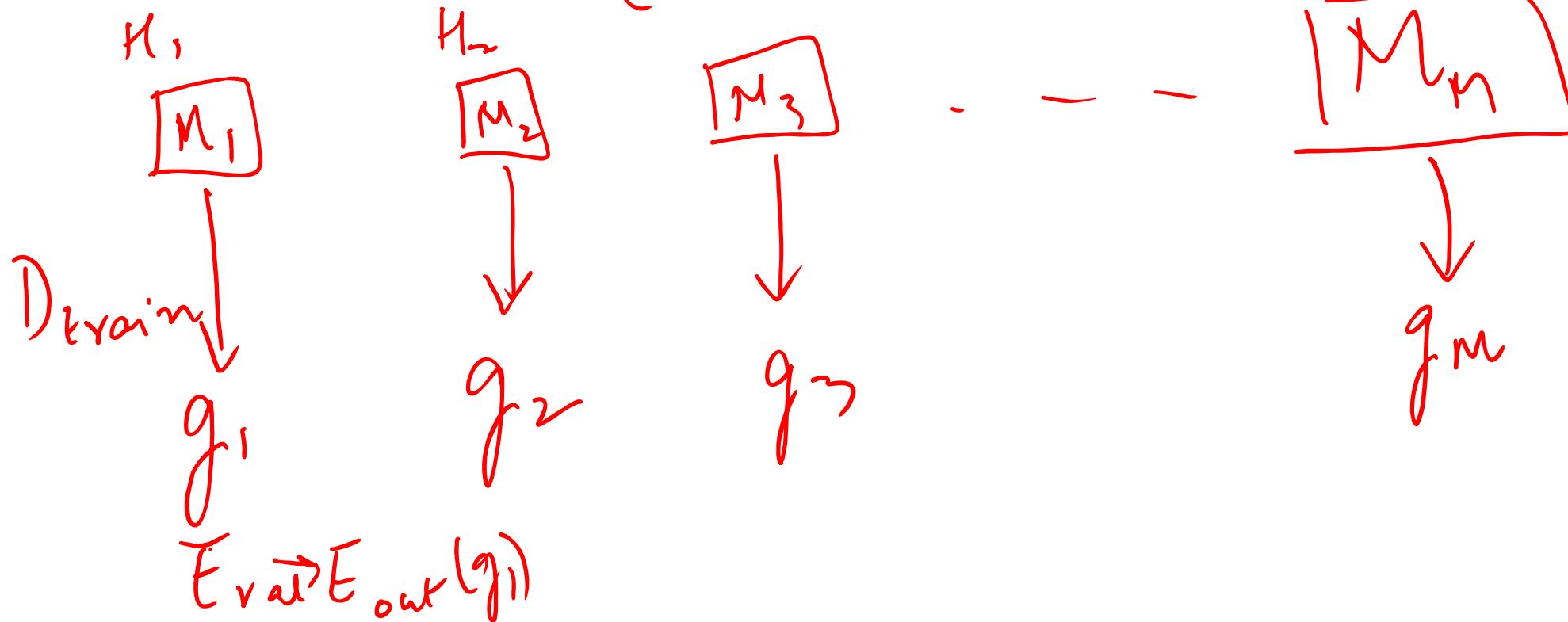
$g^-$ is behind closed doors.

$$E_{\text{out}}(g) \qquad E_{\text{out}}(g^-)$$
$$\downarrow \qquad\qquad \downarrow$$
$$E_{\text{in}}(g) \qquad E_{\text{val}}(g^-)$$

$\underbrace{\qquad\qquad\qquad\qquad}$ which should we use?

$E$

**CUSTOMER**

$g^-$

In the diagram:
$\mathcal{D}$ $(N)$
$\mathcal{D}_{\text{train}}$ $(N-K)$
$g^-$
$\mathcal{D}_{\text{val}}$ $(K)$
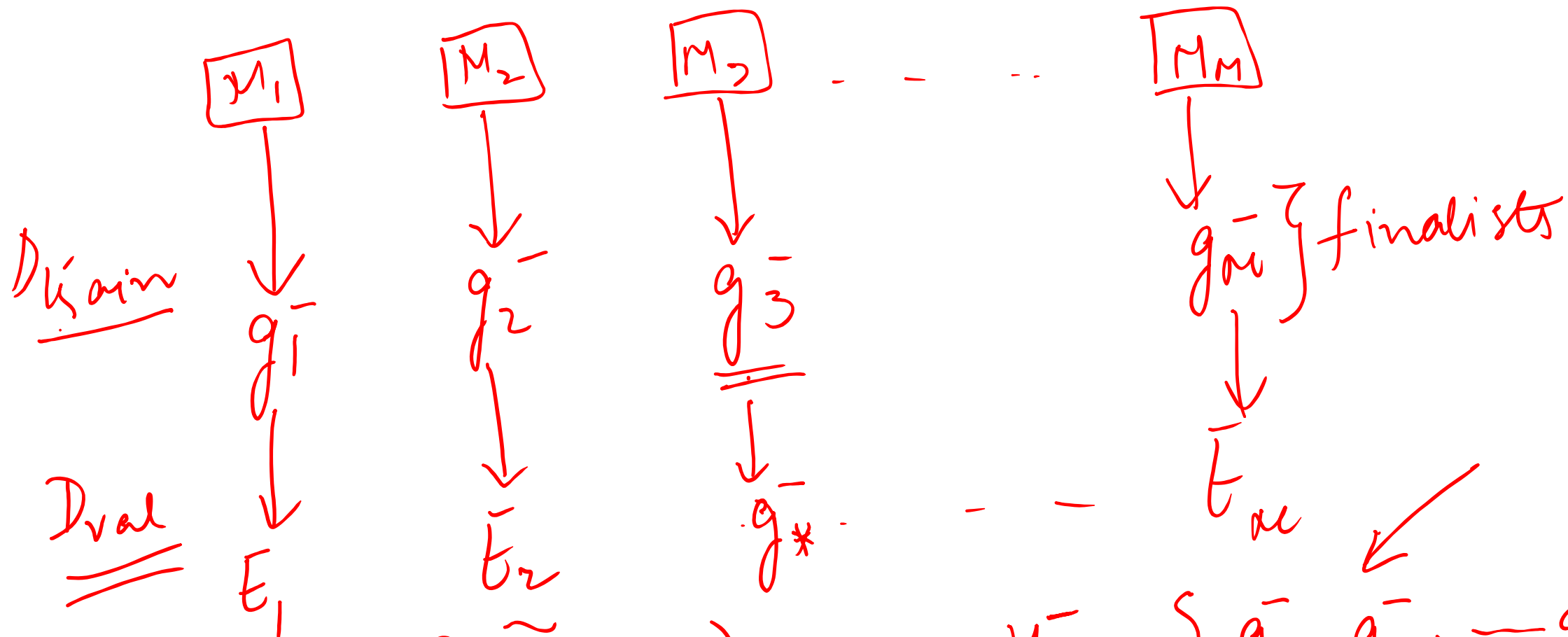$g$
$E_{\text{val}}(g^-)$

$$\underline{E_{val} \quad vs \quad E_{in}}$$

$$E_{out}(g) \leq E_{in}(g) + O\left(\sqrt{\frac{d_{vc} \ln N}{N}}\right) \leftarrow \text{H-set}$$

$$E_{out}(g) \leq E_{out}(g^-) \leq E_{val} + O\left(\frac{1}{\sqrt{K}}\right) \rightarrow$$

$$\uparrow g^- \qquad \qquad \uparrow \qquad \rightarrow (\sqrt{K}) \\ \text{single hypothesis}$$

# Model Selection

- $D_N \rightarrow D_{N-K}$ (train) $\underline{D_{train}}$

  $D_K$ (validation) $D_{val}$

$H_1$

$\boxed{M_1}$

$H_2$

$\boxed{M_2}$

$\boxed{M_3}$ $- - -$ $\boxed{M_n}$

$D_{train}$

$g_1$

$g_2$

$g_3$

$g_M$

$E_{val} \vec{\ } E_{out}(g_1)$

$$\boxed{M_1} \qquad \boxed{M_2} \qquad \boxed{M_3} \quad - \quad - \quad - \quad \boxed{M_M}$$

$D_{train}$

$$\bar{g_1} \qquad \bar{g_2} \qquad \underline{\bar{g_3}} \qquad\qquad \bar{g_M} \} finalists$$

$D_{val}$

$$E_1 \qquad \bar{E_2} \qquad \bar{g}* \quad - \quad - \quad - \quad \bar{E}_{M}$$

$$\simeq E_{out}(\bar{g_1}) \simeq E_{out}(\bar{g_2}) \qquad\qquad\qquad H^- = \{\bar{g_1}, \bar{g_2} \cdots \bar{g_M}\}$$

$$\text{finite hyb.} \qquad E_{out}(\bar{g}*) \leq E_{val}(\bar{g}*) + O\left(\sqrt{\frac{\ln M}{K}}\right)$$

$$\text{sets} \qquad\qquad\qquad\qquad \uparrow \text{finite}$$

$$E_{out}(g_*) \leq E_{out}(g^-_{m*}) \leq E_{val}(g^-_{m*}) + O\left(\sqrt{\frac{\ln M}{K}}\right)$$

$\lambda =$

$\lambda_1 = 0$     $\lambda_2 = 0.001$        $\lambda_M = 0.1$

$\boxed{\mathcal{H}_1}$    $\boxed{\mathcal{H}_2}$   - - - - -   $\boxed{\mathcal{H}_M}$

$g^-$

$D_{train}$  → $g^-_1$

$D_{val}$  ↓      ↓           ↓   $g$

$E_1$     $E_2$          $E_M$

# Validation (general philosophy)

$$g \approx g^-$$

$$K, \quad N-K$$

$$g^-(g)$$

$$E_{out}$$

$$E_{out}(g) \simeq E_{out}(g^-) \simeq E_{val}(g^-) + O\left(\frac{1}{\sqrt{K}}\right)$$

K small $\rightarrow$ tight

K big $\rightarrow$ tight

have

Can we $\sim$ $K = 1$ ?

# CROSS VALIDATION. ($K=1$?)



$$E(e_1) = E_{out}(g_1^-)$$

$$E(e_2) = E_{out}(g_2^-)$$

$$E(e_3) = E_{out}(g_3^-)$$

Pick all of them $\longrightarrow$ An average

$$E_{CV} = \frac{1}{N} \sum_{n=1}^{N} e_n \longrightarrow$$

leave one out

$(x_n, y_n)$

$$E_{\mathrm{cv}} = \frac{1}{N} \sum_{n=1}^{N} \mathrm{e}_n$$

$$\underline{\text{Theorem}} \qquad E_{D_N} \underbrace{[\overline{E}_{cv}]}_{\xi} = E_{out}(N-1)$$

$$\underbrace{E_{out}(g) \leq \overline{E}_{cv}}_{} \quad \overset{\xi \cdot}{}$$

$$\overline{E}_{out} \leq E_{cv} + O\left(\frac{1}{\sqrt{N}}\right) \to \text{average}$$

$$N's \longrightarrow e_1 \ e_2 \ e_3 \cdots e_N$$

Assume almost independent.

$$e_1 \longrightarrow (x_1 y_1) \to \overline{g}_1$$

$$e_2 \longrightarrow (x_2 y_2) \to \overline{g}_2$$

subtle way

# K-fold cross validation.

N - data points $\longrightarrow$ N regression problems

$\therefore$ N+1 training sets.

10 fold CV $\rightarrow$ 10% of your data



$\bar{g_1}$ | $\bar{g_2}$ |

10%

Analytically $\longrightarrow$ Regression

$E_{cv}(\lambda_1) \quad E_{cv}(\lambda_2) \ldots$

$\rightarrow$ Digits data

① Get features : Symmetry & Intensity.

② PLA

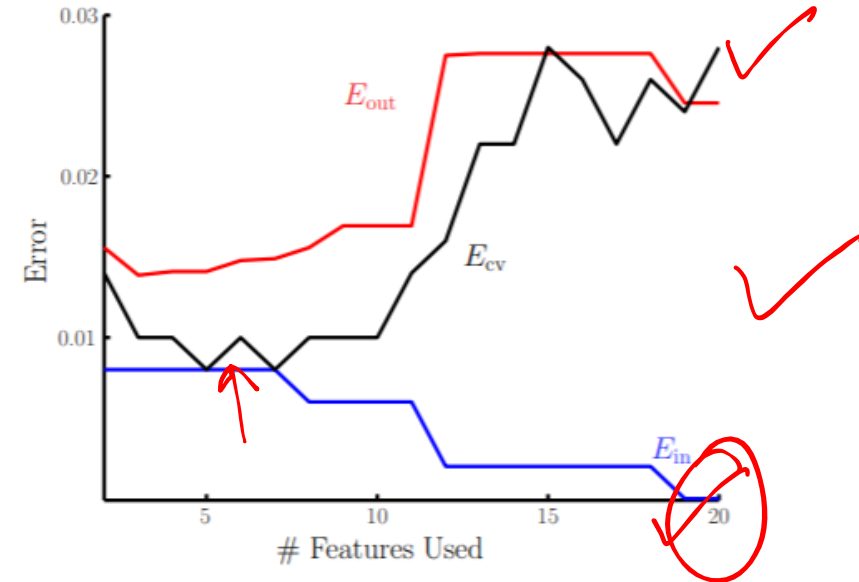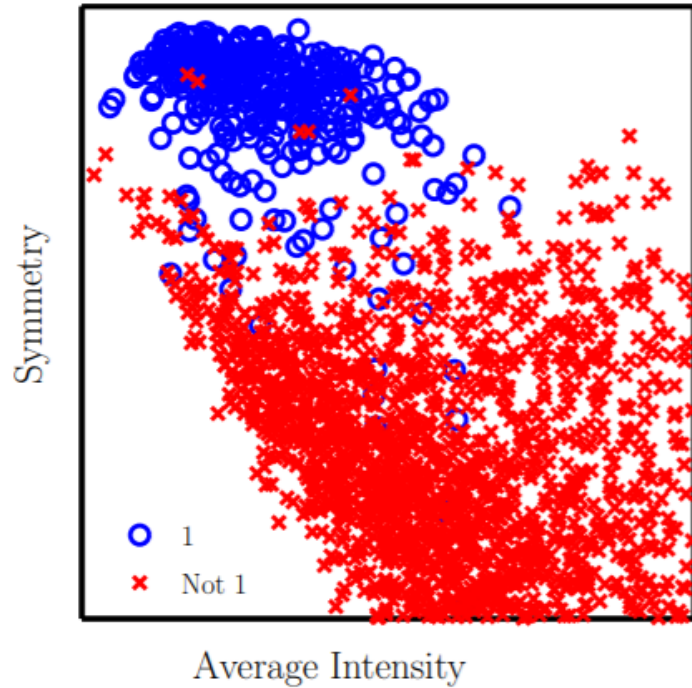③ $\begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \xrightarrow{\ \Phi\ } \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1^5 \\ x_1^5 \\ x_2^5 \end{bmatrix}$ 20 dim in $Z$ space ✓

④ Model selection

$M_1 \rightarrow \begin{bmatrix} 1 \\ z_1 \end{bmatrix}$    $M_2 \rightarrow \begin{bmatrix} 1 \\ z_1 \\ z_2 \end{bmatrix}$  - - - $M_{20}$
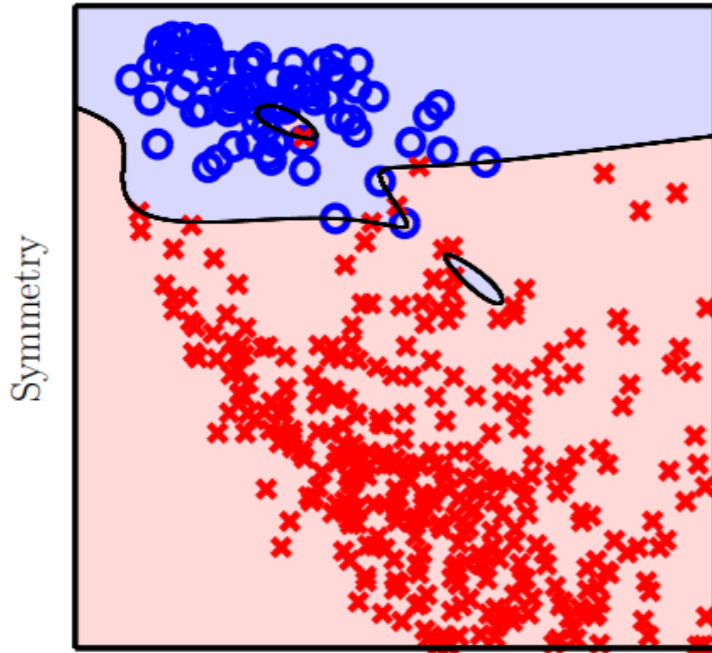
⑤ CV to pick $M_i$

# Digits Problem: '1' Versus 'Not 1'

2 cases



Symmetry

○ 1
✗ Not 1

Average Intensity

Error — $E_{\text{out}}$, $E_{\text{cv}}$, $E_{\text{in}}$

\# Features Used

$$\mathbf{x} = (1, x_1, x_2)$$

$$\mathbf{z} = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2, x_1^3, x_1^2 x_2, x_1 x_2^2, x_2^3, \ldots, x_1^5, x_1^4 x_2, x_1^3 x_2^2, x_1^2 x_2^3, x_1 x_2^4, x_2^5)$$

5th order polynomial transform $\longrightarrow$ 20 dimensional non linear feature space

# Validation Wins In the Real World



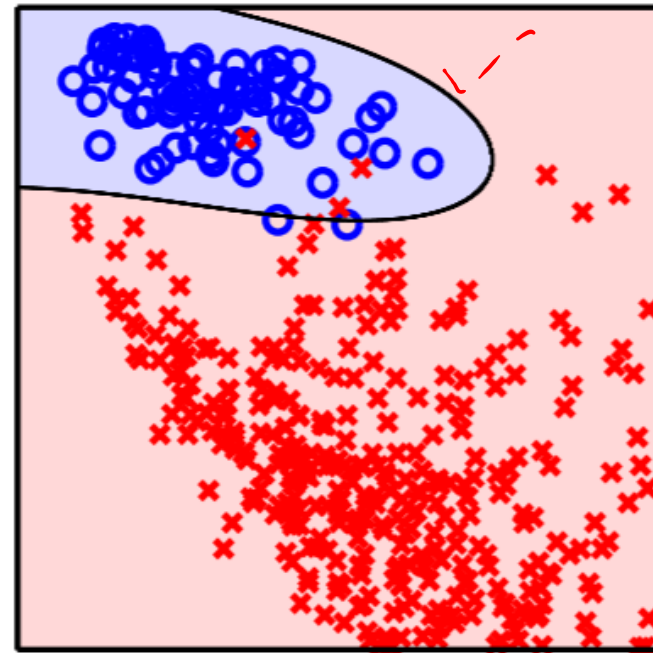no validation (20 features)

$E_{\text{in}} = 0\%$
$E_{\text{out}} = 2.5\%$

cross validation (6 features)

$E_{\text{in}} = 0.8\%$
$E_{\text{out}} = 1.5\%$

# Thanks!