# 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.
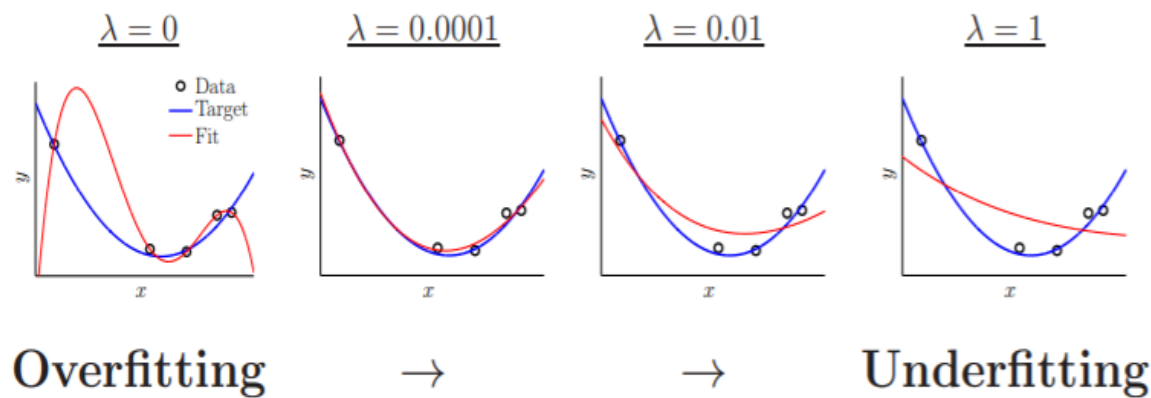
$$E_{\text{aug}}(h) = E_{\text{in}}(h) + \frac{\lambda}{N}\Omega(h)$$

$\Omega(h) \rightarrow$ 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.

# Validation

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

$\underbrace{\phantom{\text{overfit penalty}}}$

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

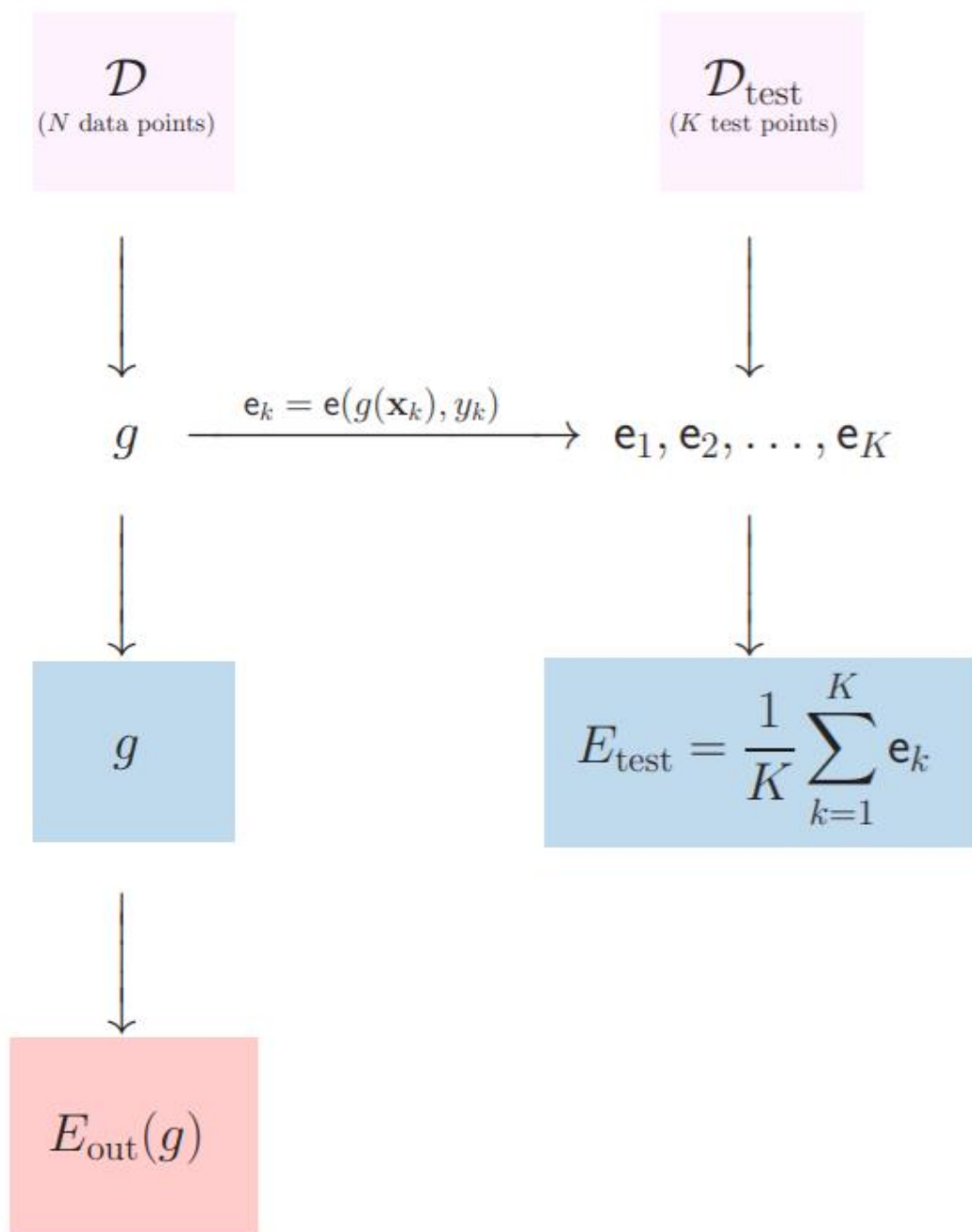regularization estimates this through a heuristic complexity penalty $\Omega(g)$

Validation goes directly for the jugular:

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

$\underbrace{\phantom{\text{overfit}}}$

validation estimates this directly

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

$\mathcal{D}$

($N$ data points)

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

($K$ test points)

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

$g$

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

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

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

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

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

$$= \frac{1}{K} \sum_{k=1}^{K} E_{\text{out}}(g) = E_{\text{out}}(g)$$
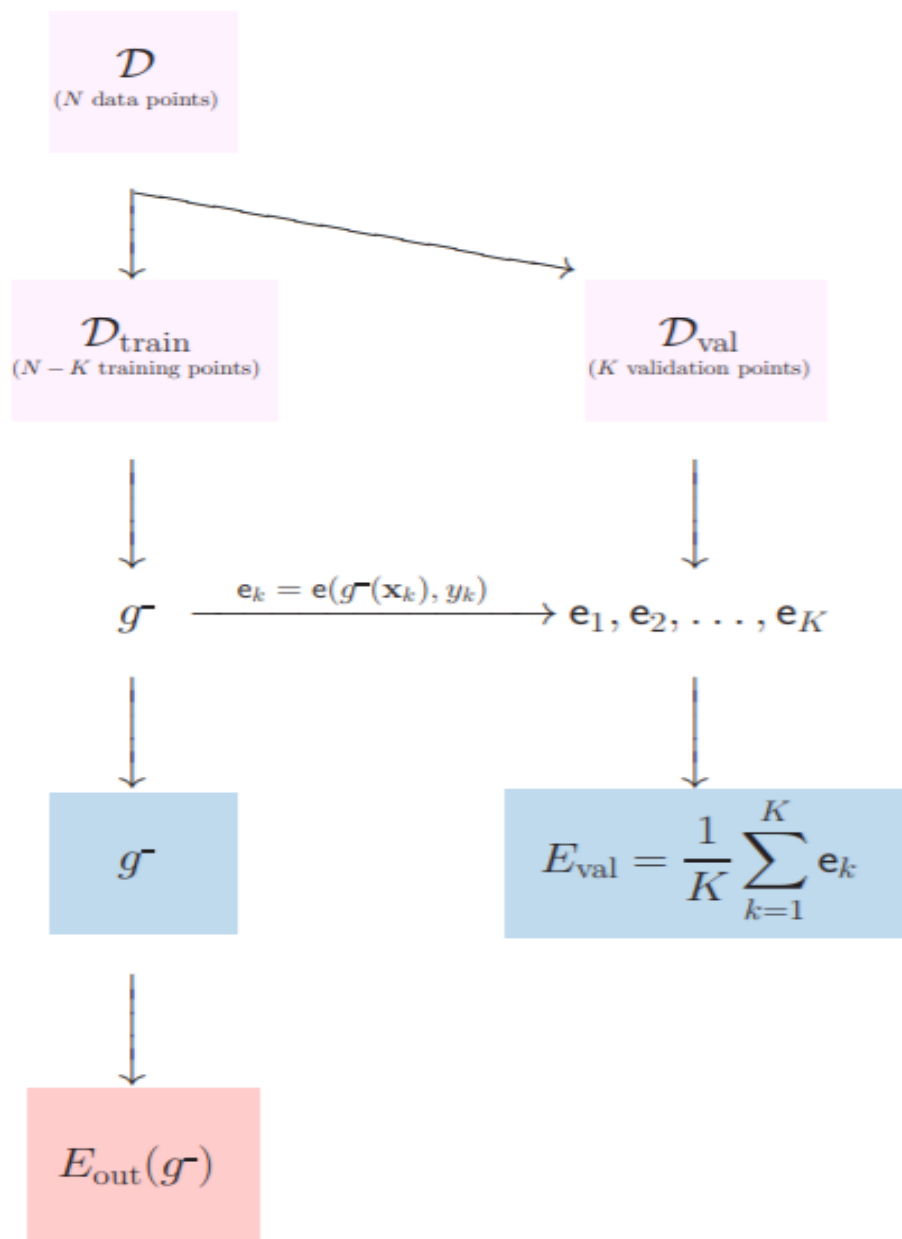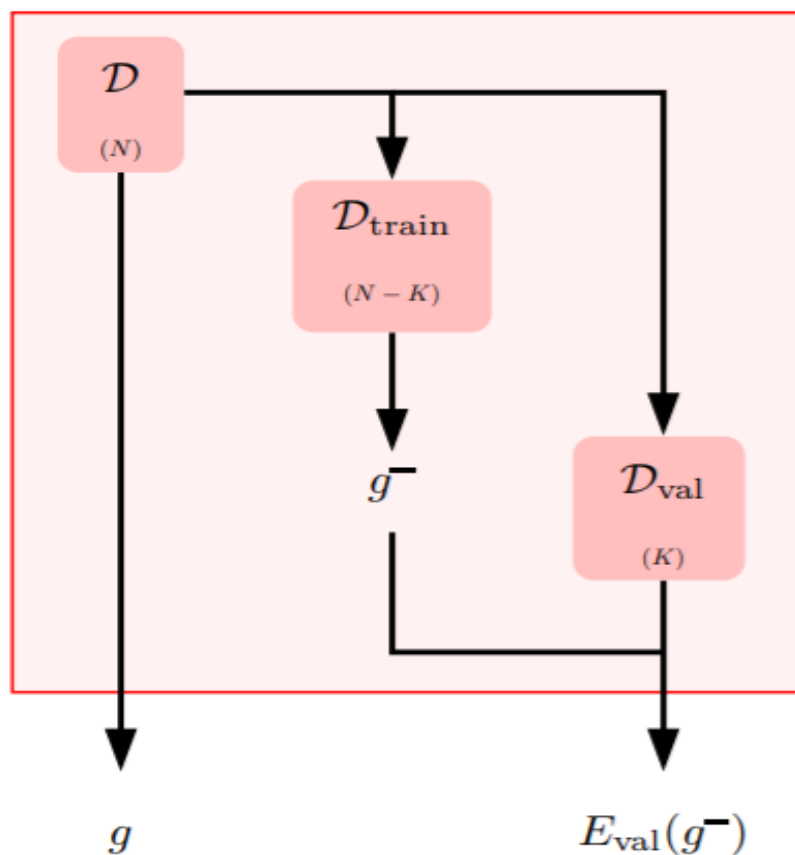
$\mathbf{e}_1, \ldots, \mathbf{e}_K$ are *independent*

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

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

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

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

# The Validation Set

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

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

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

$g^- \xrightarrow{\quad e_k = e(g^-(\mathbf{x}_k), y_k) \quad} e_1, e_2, \ldots, e_K$

$g^-$

$E_{\text{val}} = \dfrac{1}{K} \sum_{k=1}^{K} e_k$

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

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

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

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

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

$e_1, \ldots, e_K$ are *independent*

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

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

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

# Restoring $\mathcal{D}$



**Primary goal:** output best hypothesis.

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

**Secondary goal:** estimate $E_{\text{out}}(g)$.

$g^-$ is behind closed doors.

$$
\begin{array}{cc}
E_{\text{out}}(g) & E_{\text{out}}(g^-) \\
\downarrow & \downarrow \\
E_{\text{in}}(g) & E_{\text{val}}(g^-)
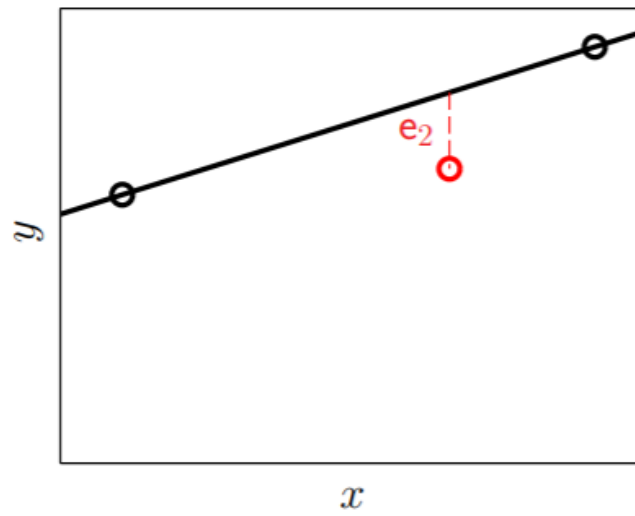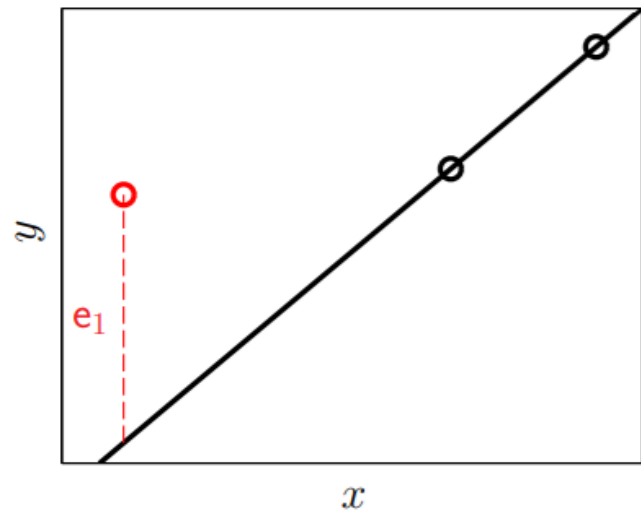\end{array}
$$

which should we use?
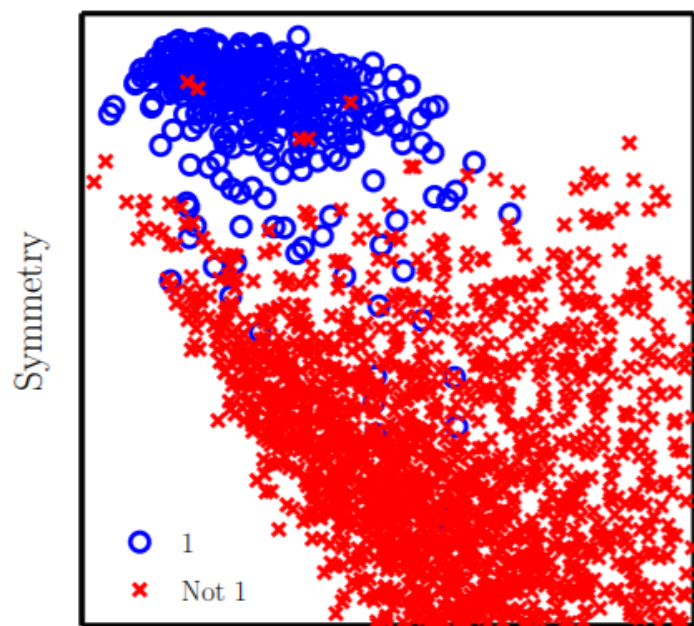
# Model Selection

-

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

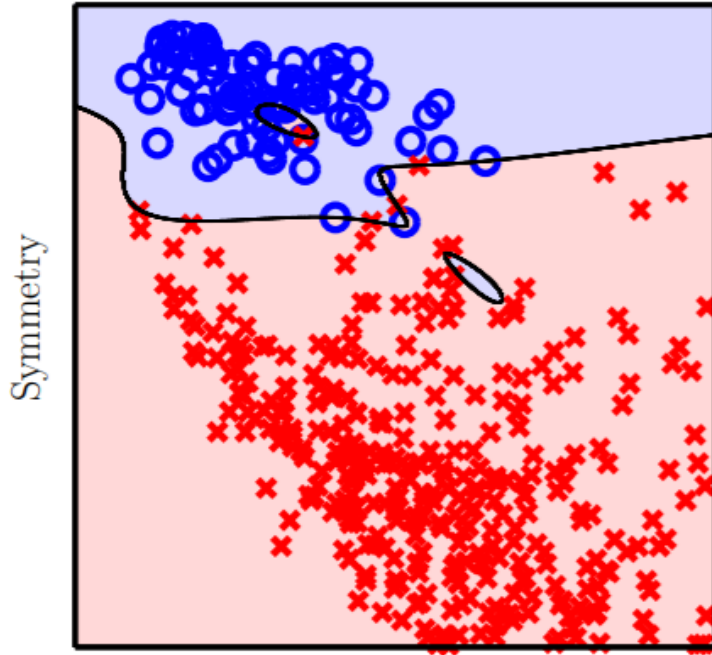# Digits Problem: '1' Versus 'Not 1'



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

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

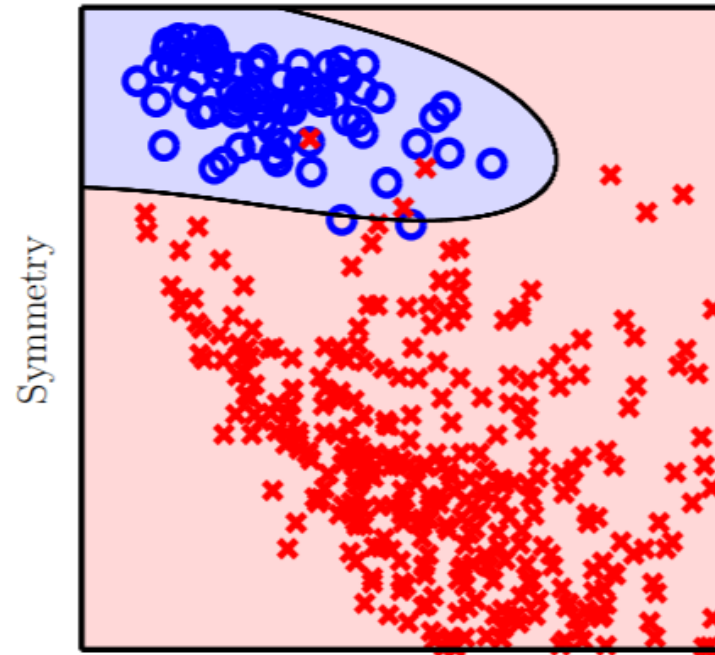5th order polynomial transform ⟶ 20 dimensional non linear feature space

# Validation Wins In the Real World



no validation (20 features)

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

cross validation (6 features)

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

# Thanks!