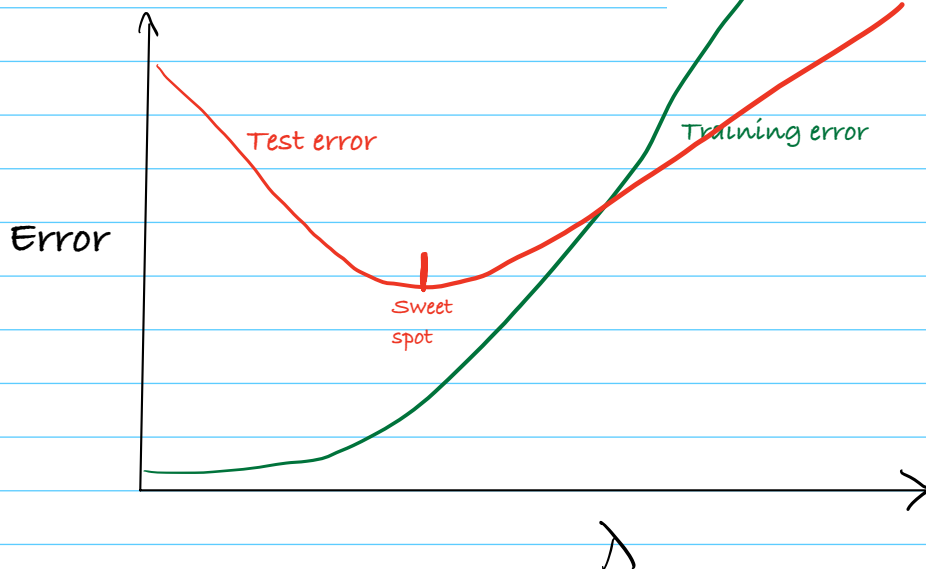# Quiz:

1. Which of the following is a better proxy for generalization error?
   A. Average training loss
   B. Average validation set loss

2. When n is large and hypothesis set size k is small, training error is close to generalization error with high probability (of $1 - \delta$), what is this probability over?

How to choose parameters?

Eg: $\min_{\mathbf{w}} \dfrac{1}{n} \sum_{i=1}^{n} \underbrace{\ell(h_{\mathbf{w}}(\mathbf{x}_i), y_i)}_{Loss} + \underbrace{\lambda r(w)}_{Regularizer}$



Idea: Use validation set to pick parameters

(pick one with smallest validation error)

## Validation Set

| Training | Validation |
|---|---|

1. If validation set is too small we dont get good estimate of error
2. But if it is very big then training set size is small

## K- fold cross validation

Partition data into K folds $D_1, ..., D_K$

For $\lambda \in \{0.01, 0.02, ..., 1\}$

$D_{-k} = D \setminus D_k$

    For k = 1 to K:

$$h_k = \text{Algorithm}(D_{-k}, \lambda) \quad \text{\% train on all data But } D_k$$

$$\varepsilon_{k,\lambda} = \frac{1}{|D_k|} \sum_{(x,y)\in D_k} \ell(h(x), y) \quad \text{\% evaluate on } D_k$$

$$\varepsilon_\lambda = \frac{1}{K} \sum_{k=1}^{K} \varepsilon_{k,\lambda} \quad \text{\% Take average validation error}$$

$$\lambda^* = \arg\min_\lambda \varepsilon_\lambda \quad \text{\% Pick best } \lambda$$



$D$

$k = 1 \; t=5$

Eg

Data



$R = 2$

$R = >$

validate on $D_3$

Train on $\{D_1, D_2, D_4, D_5\}$

Search for parameters

1. Grid search 2. Random search

3. Zooming in

# Bias Variance Decomposition

Consider the regression problem where given x we want to predict a real valued outcome y.

$$\ell(h(x), y) = (h(x) - y)^2$$

Eg

$$x = (\#bedrooms, \ sq \ footage \ of \ house, \ \# bathrooms, \dots)$$

$$y = house \ price$$

$$\text{Test error}(h) = \mathop{E}_{(x,y) \sim P} \left[(h(x) - y)^2\right]$$

Bayes Optimal Predictor = hypothesis with smallest possible test loss

$$\bar{y}(x) = E[y \mid x]$$   is the Bayes optimal predictor

$$\mathop{E}_{(x,y) \sim P} \left[(\bar{y}(x) - y)^2\right] = \text{inherent (unavoidable) noise}$$

No method with any amount of data can beat the above test loss

Say we have an Algorithm that takes as input dataset D and outputs hypothesis $h_D$. But $h_D$ depends on sample D which is a random draw. We are interested in understanding the expected test error

$$\text{Expected Test Error} = \mathop{E}_{D \sim P^n} \left[ \mathop{E}_{(x,y) \sim P} \left[ (h_D(x) - y)^2 \right] \right]$$

Since $h_D$ is a random let us consider its expected behavior:

$$\bar{h} := \mathop{E}_{n \sim P^n}[h_D] \qquad ie \qquad \bar{h}(x) = \mathop{E}_{D \sim P^n}[h_D(x)]$$

Bias : $\displaystyle \mathop{E}_{x \sim P} \left[ \left( \bar{h}(x) - \bar{y}(x) \right)^2 \right]$

<span style="color:green">Expected sq. distance between expected model of our Algorithm and the best possible model</span>

$\displaystyle \mathop{E}_{D} h_D = \bar{h}$

Variance : $\displaystyle \mathop{E}_{D \sim P^n} \left[ \mathop{E}_{x \sim P} \left[ \left( \bar{h}(x) - h_D(x) \right)^2 \right] \right]$

<span style="color:red">Fluctuation of Algorithm's random model around its mean</span>

We will see that:

Expected test error = <u>Bias</u> + <u>Variance</u> + Inherent noise

$\displaystyle \text{Expected Test Error} = \mathop{E}_{D \sim P^n} \left[ \mathop{E}_{(x,y) \sim P} \left[ \left( h_D(x) - y \right)^2 \right] \right]$

$\displaystyle = \mathop{E}_{D \sim P^n} \left[ \mathop{E}_{(x,y)} \left[ \underbrace{h_D(x) - \bar{y}(x)}_{A} + \underbrace{\bar{y}(x) - y}_{B} \right]^2 \right]$

$\displaystyle = \mathop{E}_{D \sim P^n} \left[ \mathop{E}_{(x,y)} \left[ \underbrace{\left( h_D(x) - \bar{y}(x) \right)^2}_{A^2} + \underbrace{2 \left( h_D(x) - \bar{y}(x) \right) \left( \bar{y}(x) - y \right)}_{2AB} + \underbrace{\left( \bar{y}(x) - y \right)^2}_{B^2} \right] \right]$

$\displaystyle = \mathop{E}_{D \sim P^n} \mathop{E}_{x,y} \left[ \left( h_D(x) - \bar{y}(x) \right)^2 \right] + 2 \underbrace{\mathop{E}_{D \sim P^n} \mathop{E}_{(x,y)} \left[ \left( h_D(x) - \bar{y}(x) \right) \left( \bar{y}(x) - y \right) \right]}_{\overset{\shortparallel}{0 \ ??}}$

$\displaystyle + \mathop{E}_{(x,y)} \left[ \left( \bar{y}(x) - y \right)^2 \right]$

$\displaystyle = \mathop{E}_{D \sim P^n} \mathop{E}_{(x,y)} \left[ \left( h_D(x) - \bar{y}(x) \right)^2 \right] + \mathop{E} \left[ \left( \bar{y}(x) - y \right)^2 \right]$
<span style="color:purple">"inherent noise"</span>

$\displaystyle = \mathop{E}_{D \sim P^n} \mathop{E}_{(x,y)} \left[ \underbrace{h_D(x) - \bar{h}(x)}_{A} + \underbrace{\bar{h}(x) - \bar{y}(x)}_{B} \right]^2 + \mathop{E} \left[ \left( \bar{y}(x) - y \right)^2 \right]$
<span style="color:purple">"inherent noise"</span>

$$= \underset{D}{E}\underset{(x,y)}{E}\left[\underbrace{(h_D(x) - \bar{h}(x))^2}_{A^2} + \underbrace{2\left(h_D(x) - \bar{h}(x)\right)\left(\bar{h}(x) - \bar{y}(x)\right)}_{2AB} + \underbrace{\left(\bar{h}(x) - \bar{y}(x)\right)^2}_{B^2}\right]$$

$$+ \; E\left[\left(\bar{y}(x) - y\right)^2\right]$$
"inherent noise"

$$= \underset{D}{E}\underset{x}{E}\left[(h_D(x) - \bar{h}(x))^2\right] + 2\underbrace{\underset{D}{E}\underset{x}{E}\left[(h_D(x) - \bar{h}(x))(\bar{h}(x) - \bar{y}(x))\right]}_{= 0 \; ??}$$
"Variance"

$$+ \; \underset{x}{E}\left[(\bar{h}(x) - \bar{y}(x))^2\right] \qquad + \; E\left[(\bar{y}(x) - y)^2\right]$$
"Bias"           "inherent noise"

$$= \underset{D}{E}\underset{x}{E}\left[(h_D(x) - \bar{h}(x))^2\right] + \underset{x}{E}\left[(\bar{h}(x) - \bar{y}(x))^2\right] \; + \; E\left[(\bar{y}(x) - y)^2\right]$$
"Variance"       "Bias"        "inherent noise"

**Variance**: Captures how much your classifier changes if you train on a different training set. How "over-specialized" is your classifier to a particular training set (overfitting)? If we have the best possible model for our training data, how far off are we from the average classifier?

**Bias**: What is the inherent error that you obtain from your classifier even with infinite training data? This is due to your classifier being "biased" to a particular kind of solution (e.g. linear classifier). In other words, bias is inherent to your model.

**Noise**: How big is the data-intrinsic noise? This error measures ambiguity due to your data distribution and feature representation. You can never beat this, it is an aspect of the data.
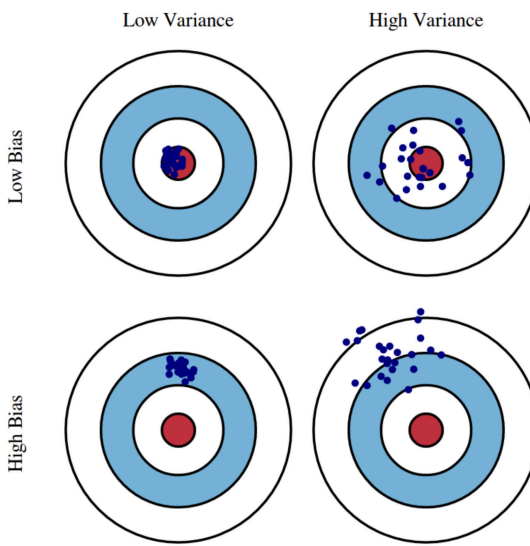
Fig 1: Graphical illustration of bias and variance. (Source http://scott.fortmann-roe.com/docs/BiasVariance.html) Fig 2: The variation of Bias and Variance with the model complexity. This is similar to the concept of overfitting and underfitting. More complex models overfit while the simplest models underfit. (Source http://scott.fortmann-roe.com/docs/BiasVariance.html)

## Detecting High Bias and High Variance

If a classifier is under-performing (e.g. if the test or training error is too high), there are several ways to improve performance. To find out which of these many techniques is the right one for the situation, the first step is to determine the root of the problem. The graph above plots the training error and the test error and can be divided into two overarching regimes. In the first regime (on the left side of the graph), training error is below the desired error threshold (denoted by $\epsilon$), but test error is significantly higher. In the second regime (on the right side of the graph), test error is remarkably close to training error, but both are above the desired tolerance of $\epsilon$.

### Regime 1 (High Variance)

In the first regime, the cause of the poor performance is high variance.

**Symptoms**:

1. Training error is much lower than test error
2. Training error is lower than $\epsilon$
3. Test error is above $\epsilon$

**Remedies**:

- Add more training data
- Reduce model complexity -- complex models are prone to high variance
- Bagging (will be covered later in the course)

### Regime 2 (High Bias)

Unlike the first regime, the second regime indicates high bias: the model being used is not robust enough to produce an accurate prediction.

**Symptoms**:

1. Training error is higher than $\epsilon$

**Remedies**:

- Use more complex model (e.g. kernelize, use non-linear models)
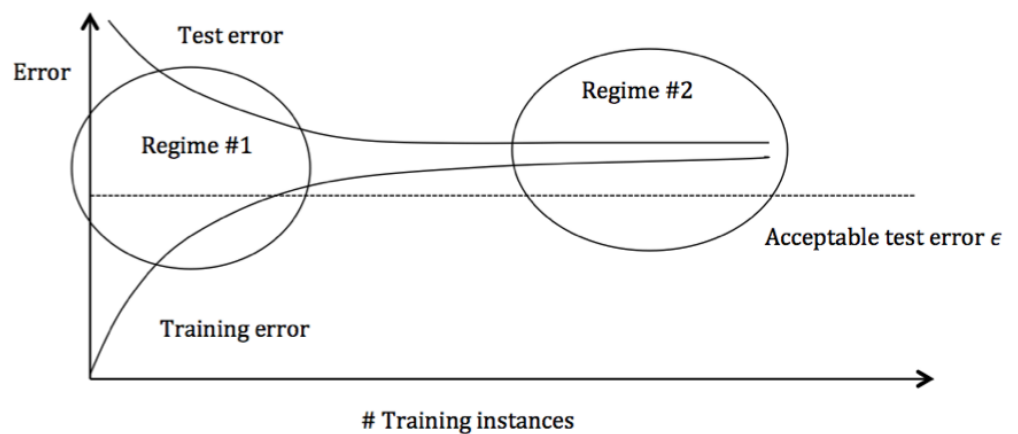- Add features



Figure 3: Test and training error as the number of training instances increases.