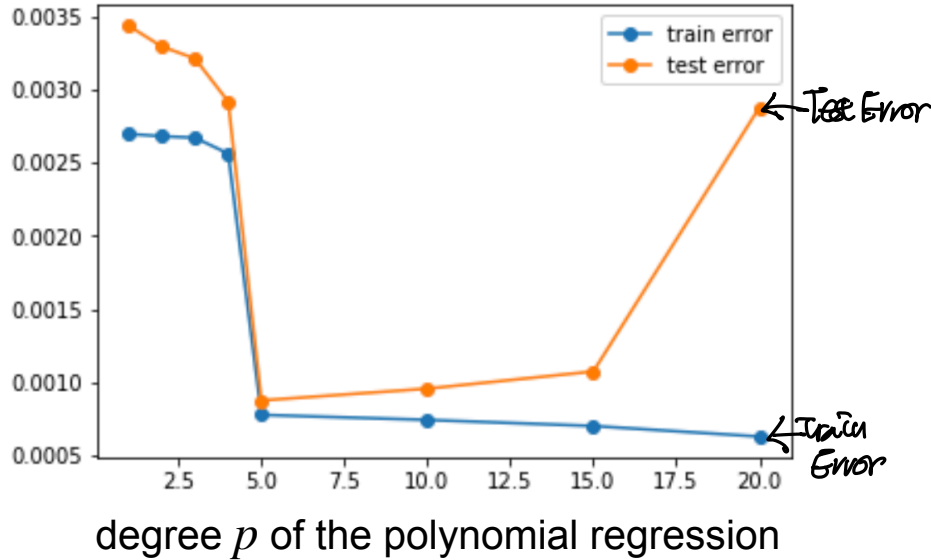# Lecture 5:
# Bias-Variance Tradeoff

- explaining test error using theoretical analysis
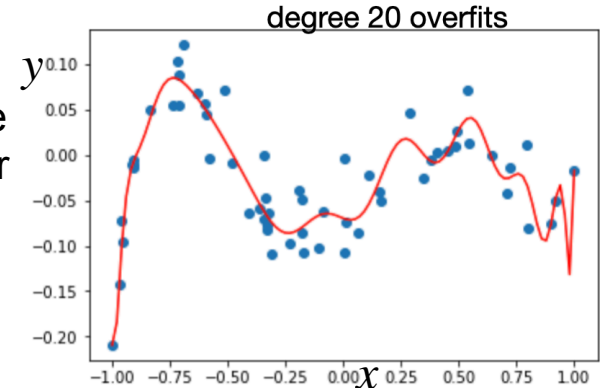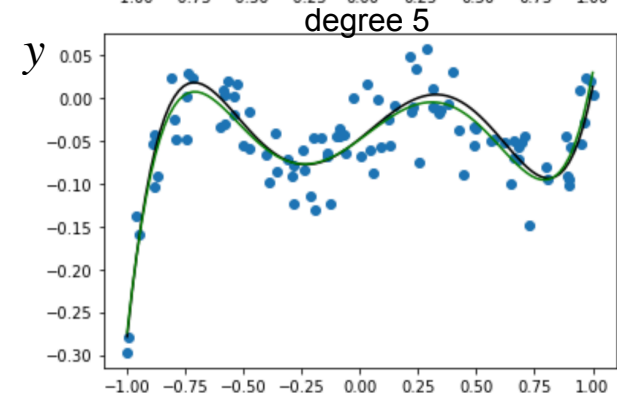
W

# Train/test error vs. complexity

Error



degree $p$ of the polynomial regression
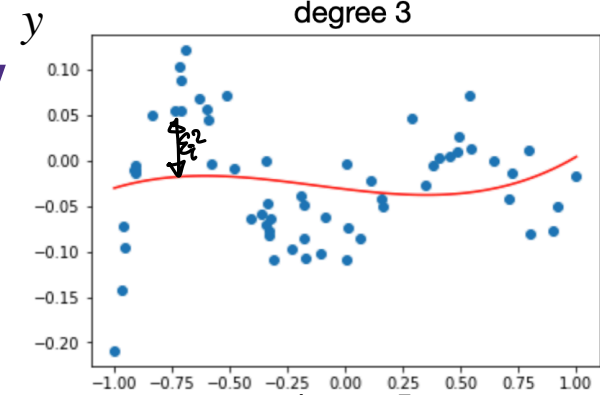
- **Model complexity** e.g., degree $p$ of the polynomial model, number of features used in diabetes example
  - Related to the dimension of the model parameter
- **Train error** monotonically decreases with model complexity
- **Test error** has a U shape

# Statistical learning

- Suppose data is generated from a statistical model $(X, Y) \sim P_{X,Y}$

  - and assume we know $P_{X,Y}$ (just for now to explain statistical learning)

- Then **learning** is to find a predictor $\eta : \mathbb{R}^d \to \mathbb{R}$ that minimizes

  - the expected error $\mathbb{E}_{(X,Y) \sim P_{X,Y}}[(Y - \eta(X))^2]$

  - think of this random $(X, Y)$ as a new sample you will encounter when you deployed your learned model, and we care about its average performance

- Since, we do not assume anything about the function $\eta(x)$, it can take any value for each $X = x$, hence the optimization can be broken into sum (or more precisely integral) of multiple objective functions, each involving a specific value $X = x$

  - $\mathbb{E}_{(X,Y) \sim P_{X,Y}}[(Y - \eta(X))^2] = \mathbb{E}_{X \sim P_X}\left[\mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 \,|\, X = x]\right]$

$$= \int \mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 \,|\, X = x] \, P_X(x) \, dx$$

Or for discrete $X$, $\quad = \sum_x P_X(x) \underbrace{\mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 \,|\, X = x]}$

*each $\eta(x)$ optimized separately*

Where we used the chain rule: $\mathbb{E}_{X,Y}[f(X, Y)] = \mathbb{E}_X\left[\mathbb{E}_{Y|X}[f(x, Y) \,|\, X = x]\right]$
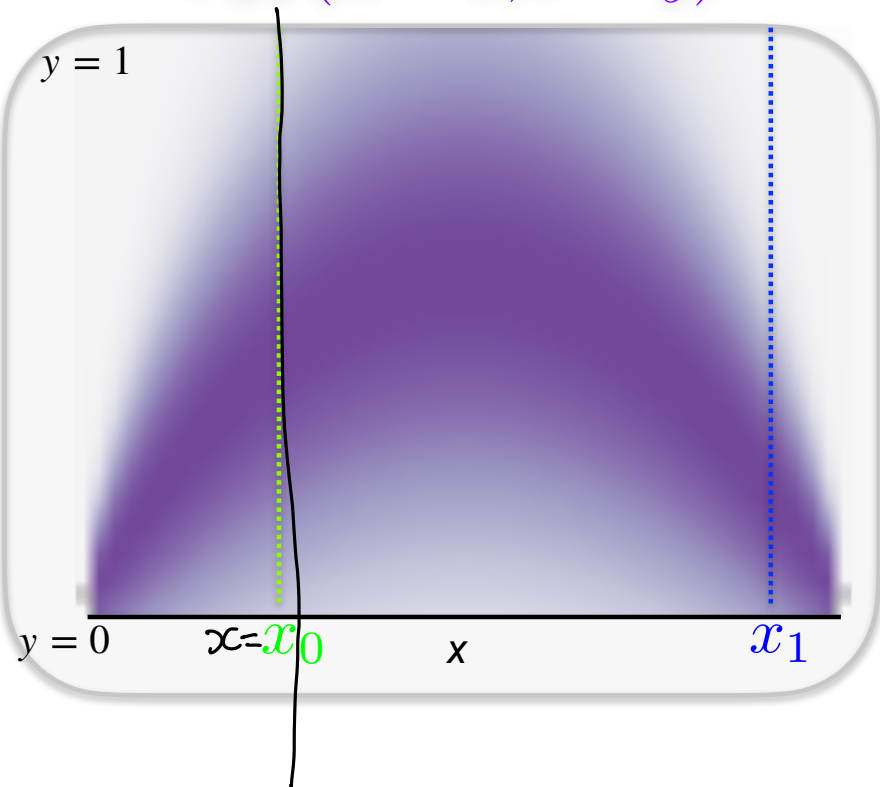
# Statistical learning

- We can solve the optimization for each *instance* $X = x$ separately

  - $\eta(x) = \arg\min\limits_{a \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}}[(Y - a)^2 \,|\, X = x]$

- The optimal solution is $\eta(x) = \mathbb{E}_{Y \sim P_{Y|X}}[Y \,|\, X = x]$,

  which is the best prediction in $\ell_2$-loss/Mean Squared Error

- Claim: $\mathbb{E}_{Y \sim P_{Y|X}}[Y \,|\, X = x] = \arg\min\limits_{a \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}}[(Y - a)^2 \,|\, X = x]$

- Proof:

  $$\frac{\partial}{\partial a} \mathbb{E}[(Y-a)^2 \,|\, X=x] = \frac{\partial}{\partial a} \left\{ \mathbb{E}[Y^2 | X=x] - 2\,\mathbb{E}[Y|X=x]\,a + a^2 \right\}$$

  $$= -2\,\mathbb{E}[Y|X=x] + 2a \Big|_{a=\eta(x)} = 0$$

  $$\mathbb{E}[Y \,|\, X=x] = \eta(x)$$

- Note that this optimal statistical estimator $\eta(x) = \mathbb{E}[Y \,|\, X = x]$ cannot be implemented as we do not know $P_{X,Y}$ in practice

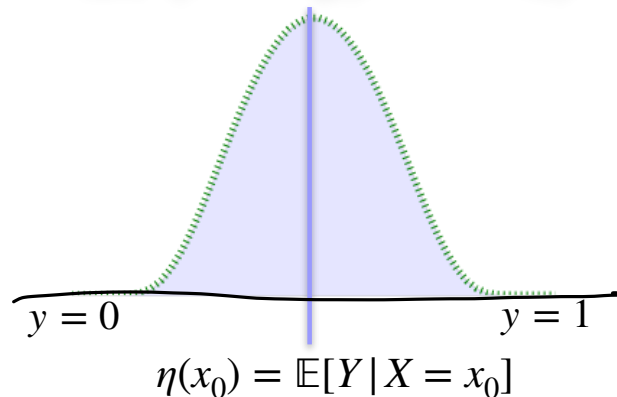- This is only for the purpose of conceptual understanding

# Statistical Learning

Ideally, we want to find:

Optimal Predictor: $\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$

$P_{XY}(X = x, Y = y)$

$y = 1$

$y = 0$    $x = x_0$    $x$    $x_1$

$P_{XY}(Y = y|X = x_0)$

$y = 0$         $y = 1$

$\eta(x_0) = \mathbb{E}[Y|X = x_0]$

$P_{XY}(Y = y|X = x_1)$

$y = 0$         $y = 1$

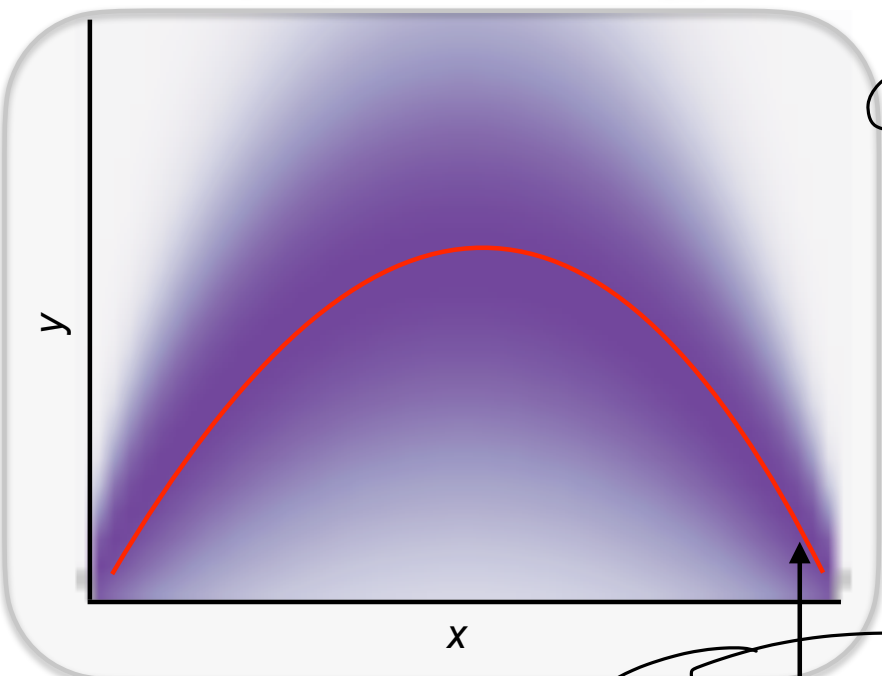$\eta(x_1) = \mathbb{E}[Y|X = x_1]$

# Statistical Learning

$$P_{XY}(X = x, Y = y)$$



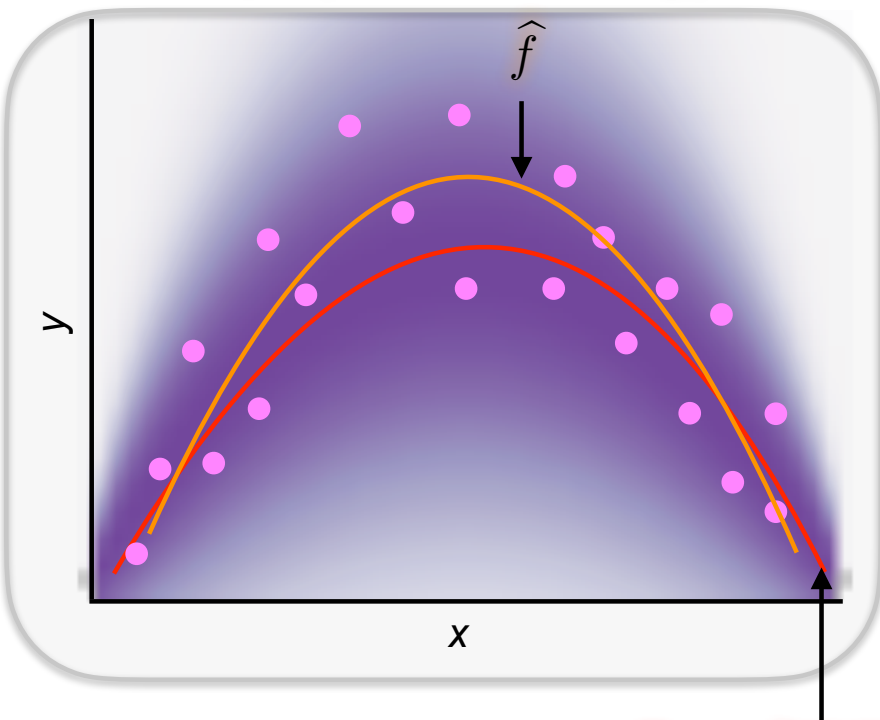Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

② But we do not know $P_{X,Y}$

We only have samples.

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

# Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
$$(x_i, y_i) \overset{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \ldots, n$$

So we need to restrict our predictor to a function class (e.g., linear, degree-$p$ polynomial) to avoid overfitting:

$$\widehat{f} = \arg\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \underbrace{(y_i - f(x_i))^2}_{\text{sample loss}}$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

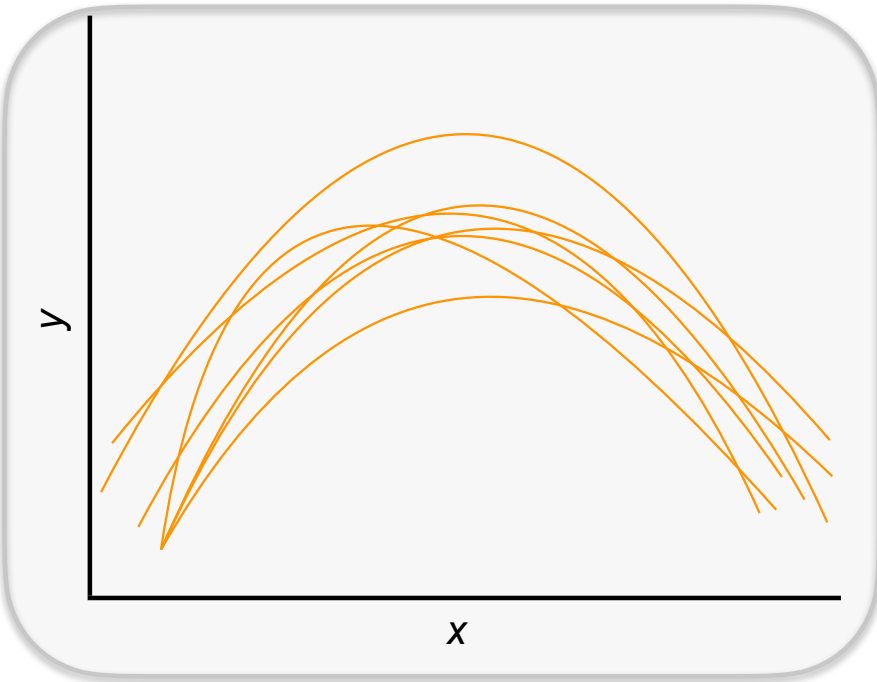We care about how our predictor performs on future unseen data
True Error of $\hat{f}$ : $\mathbb{E}_{X,Y}[(Y - \hat{f}(X))^2]$

**Future prediction error $\mathbb{E}_{X,Y}[(Y - \hat{f}(X))^2]$ is random because $\hat{f}$ is random (whose randomness comes from training data $\mathscr{D}$)**

$$P_{XY}(X = x, Y = y)$$



Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$ results in different $\widehat{f}$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- We are interested in the **True Error** of a (random) learned predictor: $\hat{f}_{\mathscr{D}}$ ;

$$\mathbb{E}_{X,Y}[(Y - \hat{f}_{\mathscr{D}}(X))^2]$$

- But the analysis can be done for each $X = x$ separately, so we analyze the **conditional true error**:

$$\mathbb{E}_{Y|X}[(Y - \hat{f}_{\mathscr{D}}(x))^2 | X = x]$$

- And we care about the **average conditional true error**, averaged over training data:

$$\mathbb{E}_{\mathscr{D}}\big[\mathbb{E}_{Y|X}[(Y - \hat{f}_{\mathscr{D}}(x))^2 | X = x]\big]$$

written compactly as $\quad = \mathbb{E}[(Y - \hat{f}_{\mathscr{D}}(x))^2]$

3 parts / 3 sources of error.

to understand this error, we decompose it into

$$\mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2] | X = x] = \mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2] | X = x]$$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathcal{D}} = \arg\min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

- **Average conditional true error**:

$$\mathbb{E}_{\mathcal{D}, Y|x}[(Y - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{\mathcal{D}, Y|x}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2]$$

$$= \mathbb{E}\left[ (Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_{0}(x)) + (\eta(x) - \hat{f}_{0}(x))^2 \right]$$

$$= \mathbb{E}_{Y|x}\left[ (Y - \eta(x))^2 \right] + 2 \cdot \mathbb{E}\left[ (Y - \eta(x))(\eta(x) - \hat{f}_{D}(x)) \right] + \mathbb{E}\left[ (\eta(x) - \hat{f}_{0}(x))^2 \right]$$

$$\to \mathbb{E}[Y - \eta(x)] \, \mathbb{E}[(\eta(x) - \hat{f}_{0}(x))]$$

$$\to \mathbb{E}[\mathbb{E}[X=x] - \eta(x)]$$

$$= 0.$$

$$= \underbrace{\mathbb{E}_{Y|x}\left[ (Y - \eta(x))^2 \right]}_{\text{Irreducible error}} + \underbrace{\mathbb{E}\left[ (\eta(x) - \hat{f}_{D}(x))^2 \right]}_{\text{Expected Learning Error.}}$$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathcal{D}} = \arg\min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

- **Average conditional true error**:

$$\mathbb{E}_{\mathcal{D}, Y|x}[(Y - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{\mathcal{D}, Y|x}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2]$$

$$= \mathbb{E}_{\mathcal{D}, Y|x}\left[(Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x)) + (\eta(x) - \hat{f}_{\mathcal{D}}(x))^2\right]$$

$$= \mathbb{E}_{Y|x}[(Y - \eta(x))^2] + 2\underbrace{\mathbb{E}_{\mathcal{D}, Y|x}[(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x))]}_{=0} + \mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]$$

(this follows from independence of $\mathcal{D}$ and $(X, Y)$ and
$\mathbb{E}_{Y|x}[Y - \eta(x)] = \mathbb{E}[Y|X = x] - \eta(x) = 0$)

$$= \mathbb{E}_{Y|x}[(Y - \eta(x))^2] \quad + \quad \mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]$$

**Irreducible error**
(a) Caused by stochastic
label noise in $P_{Y|X=x}$
(b) cannot be reduced

**Average learning error**
Caused by
*(a)* either using too "simple" of a model or
*(b)* not enough data to learn the model accurately

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i,y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

$$\mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D}}\big[ \big( \eta(x) - \overbrace{\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)]}^{\bar{f}(x)} + \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x) \big)^2 \big]$$

$$= \mathbb{E}\left[ (\eta(x) - \bar{f}(x))^2 \right] + 2\, \mathbb{E}\left[ (\eta(x) - \bar{f}(x))(\bar{f}(x) - \hat{f}_D(x)) \right] + \mathbb{E}\left[ (\bar{f}(x) - \hat{f}_D(x))^2 \right]$$

$$\longrightarrow 2\,(\eta(x) - \bar{f}(x))\,(\bar{f}(x) - \underbrace{\mathbb{E}[\hat{f}_D(x)]}_{\bar{f}(x)})$$

$$= 0$$

$$= \underbrace{\underbrace{(\eta(x) - \bar{f}(x))^2}_{\text{Bias of } \hat{f}_D(x)}}_{\text{Bias}^2} + \underbrace{\mathbb{E}_D\left[ (\bar{f}(x) - \hat{f}_D(x))^2 \right]}_{\text{Variance of } \hat{f}_D(x).}$$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f\in\mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i,y_i)\in\mathscr{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

$$\mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D}}\Big[\big(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] + \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x)\big)^2\Big]$$

$$= \mathbb{E}_{\mathscr{D}}\Big[\big(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)]\big)^2 + 2(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)])(\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x))$$

$$+ \big(\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x)\big)^2\Big]$$

$$= \underbrace{\big(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)]\big)^2}_{\textbf{biased squared}} + \underbrace{\mathbb{E}_{\mathscr{D}}\Big[\big(\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x)\big)^2\Big]}_{\textbf{variance}}$$
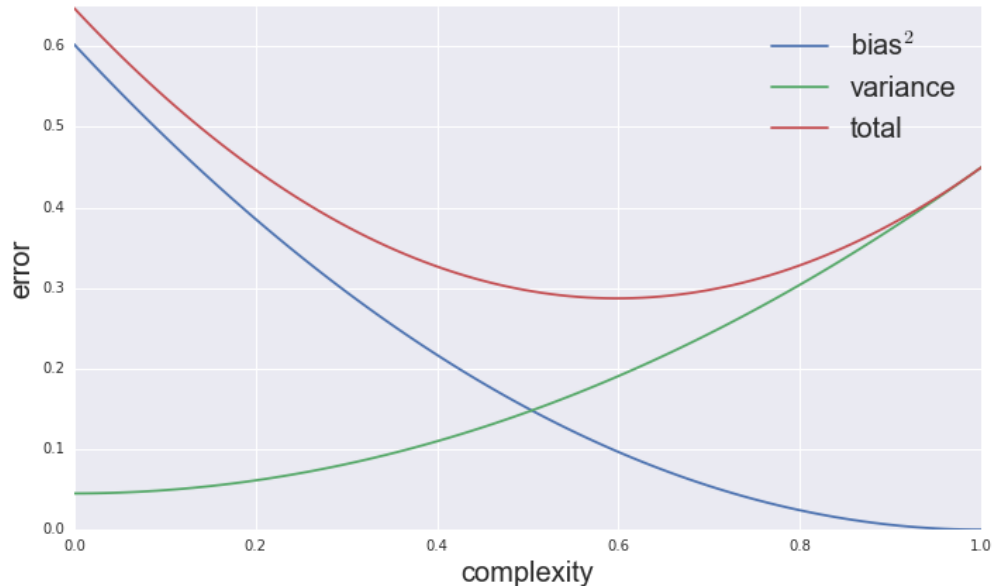
# Bias-variance tradeoff

- **Average conditional true error**:

$$\mathbb{E}_{\mathcal{D},Y|x}[(Y - \hat{f}_{\mathcal{D}}(x))^2] = \underline{\mathbb{E}_{Y|x}\left[(Y - \eta(x))^2\right]}$$

<div align="center">

**irreducible error**

</div>

$$+ \underline{\left(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)]\right)^2} + \underline{\mathbb{E}_{\mathcal{D}}\left[\left(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x)\right)^2\right]}$$

<div align="center">

**biased squared**          **variance**

</div>

**Bias squared:**
measures how the predictor is mismatched with the best predictor in expectation

**variance:**
measures how the predictor varies each time with a new training datasets

# Questions?