# Announcements

## Assignments
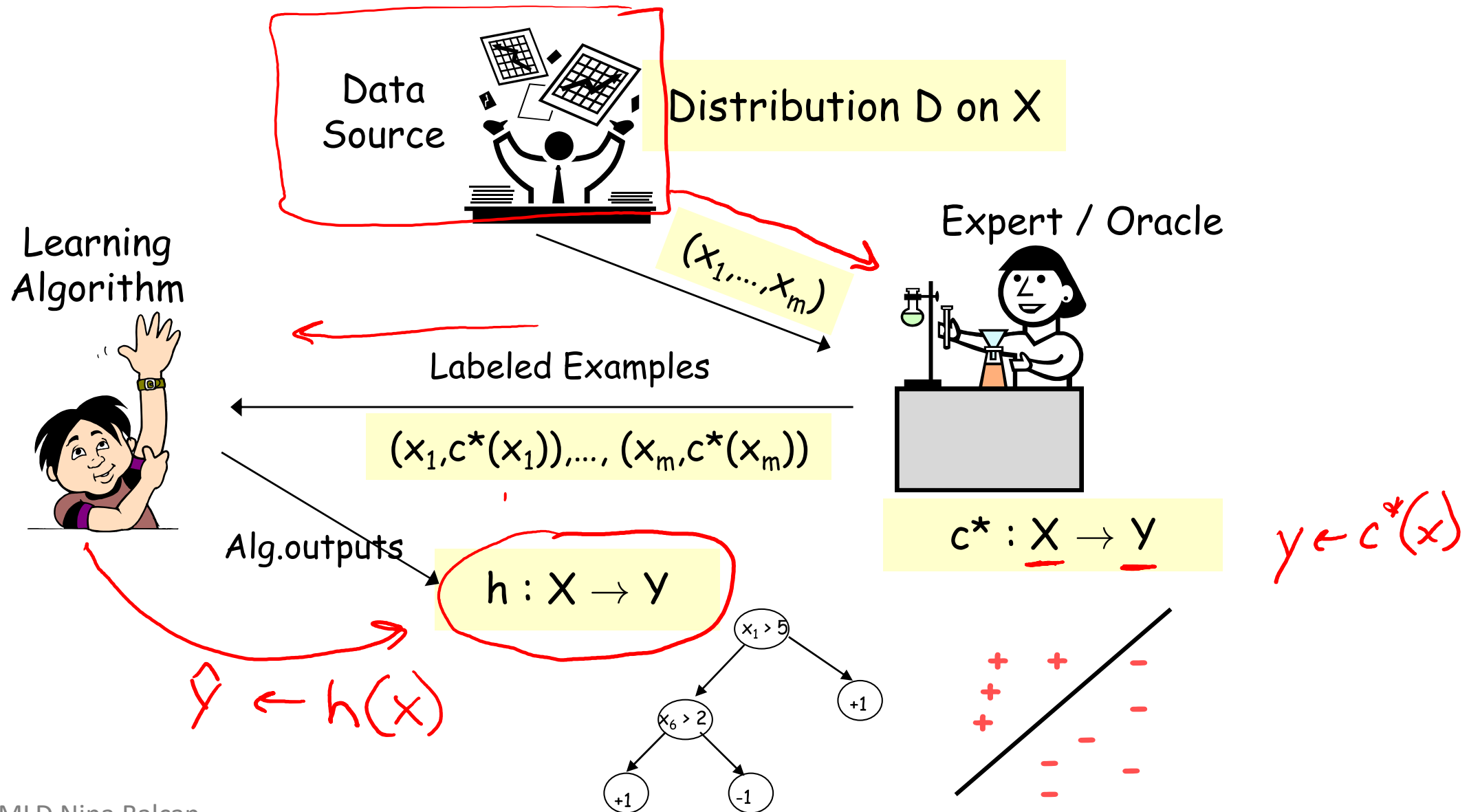
- HW10 (programming + "written")
  - Due Thu 4/30, 11:59 pm

# Introduction to Machine Learning

## Learning Theory

Instructor: Pat Virtue

# Model for Supervised Learning

# Optimal Classification Function

Find the best $h(x) \rightarrow \hat{y}$ by searching in the space of hypothesis functions $h \in \mathcal{H}$.

Optimal classifier:

labels input

$$h^*(x) = \text{argmax}_y P(Y = y \mid X = x)$$

But why?

$\arg\max_\theta \; g(x^T \theta)$

$g \, f$

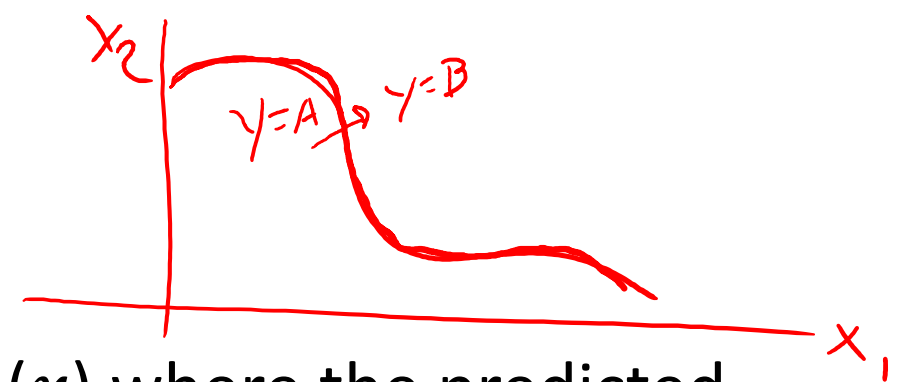$P(x \mid y) P(y)$

# Optimal Decision Boundaries

## Decision boundary

- The set of points in the domain of the input ($x$) where the predicted classification changes

## Two class decision boundary

- So far, we have decided to let the decision boundary be all $x$ such that:

$$p(Y = 0 \mid X = x) = p(Y = 1 \mid X = x)$$

- What assumptions are we making here?
  - This assumes that the cost of predicting it wrong is the same for both classes

# Optimal Classification Function

Find the best $h(x) \to \hat{y}$ by searching in the space of hypothesis functions $h \in \mathcal{H}$.

Optimal classifier:

$$h^*(x) = \underset{y}{\operatorname{argmax}} \, P(Y = y \mid X = x)$$

But why?

Goal: find a prediction function $h^* : \mathcal{X} \to \mathcal{Y}$ that minimizes the expected loss for randomly drawn test data $(X, Y)$

$$h^* = \underset{h}{\operatorname{argmin}} \, \mathbb{E}_{XY}\left[L(Y, h(X))\right]$$

$L(y, \hat{y})$ is the loss or cost of predicting $\hat{y}$ when the true value is $y$.

# Loss Functions

$$R(h)$$

$$h^* = \operatorname*{argmin}_{h} \mathbb{E}_{XY}\big[L(Y, h(X))\big]$$

Loss function:

$$L: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$$

Classification:

- Two-class, 0,1 loss

pred

|  | $\hat{y}=0$ | $\hat{y}=1$ |
|---|---|---|
| $y=0$ | 0 | 1 |
| $y=1$ | 1 | 0 |

True

- Two-class, arbitrary loss

|  | $\hat{y}=0$ | $\hat{y}=1$ |
|---|---|---|
| $y=0$ | 0 | 1 |
| SPAM $y=1$ | 100 | 0 |

NOT

pred

|  | $\hat{y}=0$ | $\hat{y}=1$ |
|---|---|---|
| $y=0$ | TN | FP |
| $y=1$ | FN | TP |

False positives and false negatives:

# Loss Functions

$$h^* = \operatorname*{argmin}_h \mathbb{E}_{XY}\big[L\big(Y, h(X)\big)\big]$$

Loss function:

$$L: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$$

Classification:

- Two-class, 0,1 loss

- Two-class, arbitrary loss

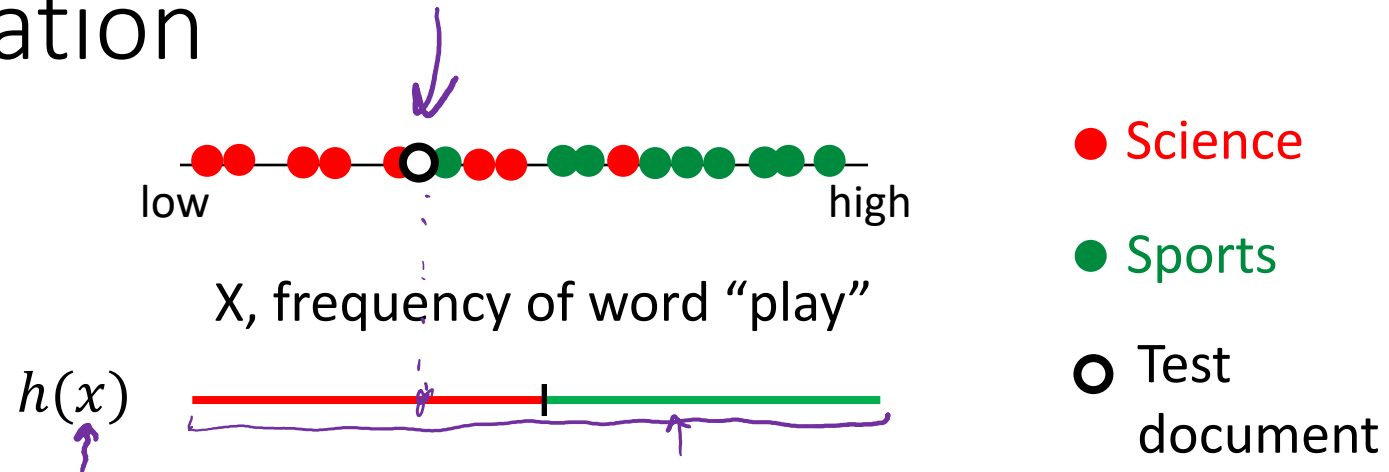Regression: $$L\big(y, h(x)\big) = \big(y - h(x)\big)^2$$

# Expected Value

Quick review

$$g(Z)$$
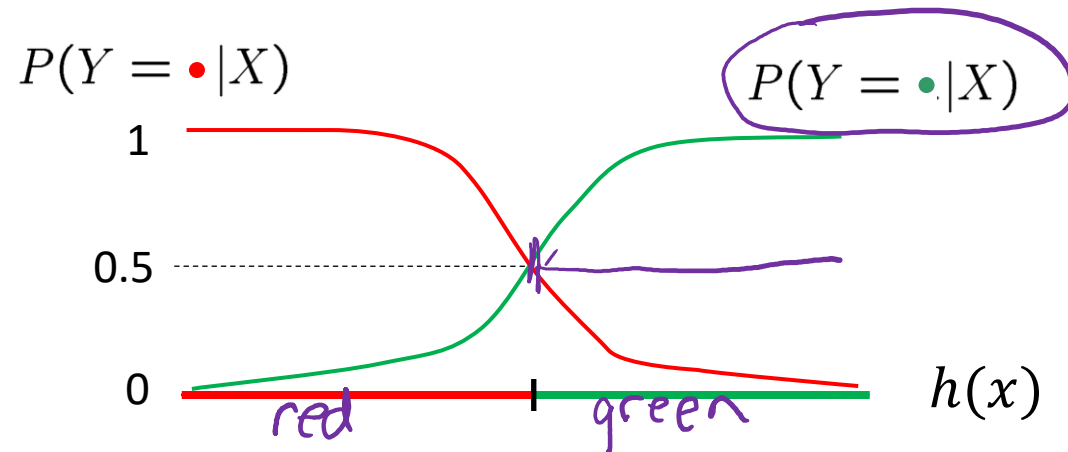
$$E[g(Z)] = \sum_{z \in Z} p(z) g(z) \quad \longleftarrow$$

$$E[g(Z)] = \int_{\mathbb{R}} f(z) g(z) \, dz$$

# Binary Classification

low          X, frequency of word "play"          high

● Science

● Sports

○ Test document

$h(x)$

Model X and Y as random variables

$P(Y = \bullet|X)$          $P(Y = \bullet|X)$

1

0.5

0

red          green          $h(x)$

For a given $x$, $h(x)$ = label Y which is more likely
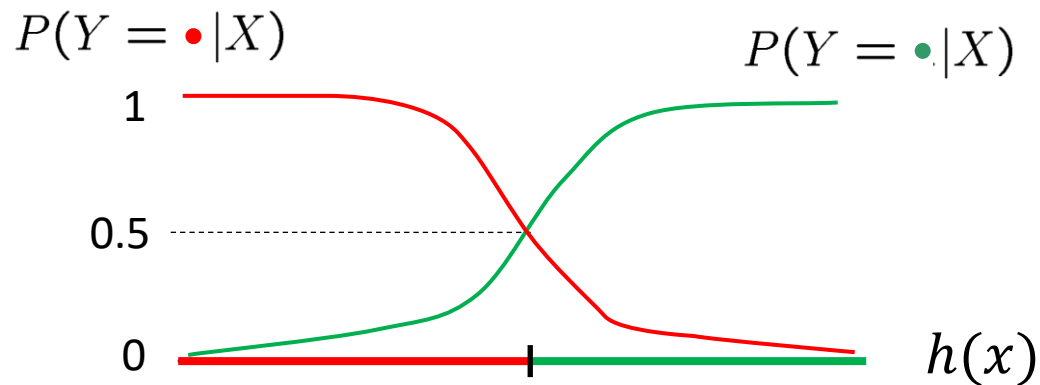
$$h(x) = \arg\max_{Y=y} P(Y = y|X = x)$$

# Optimal Classification Function

$$h^*(x) = \operatorname*{argmax}_{y \in \{0,1\}} P(Y = y \mid X = x)$$

$$h^* = \operatorname*{argmin}_{h \in \mathcal{H}} \mathbb{E}_{XY}\big[L(Y, h(X))\big]$$

Start with arbitrary two-class loss $L(y, \hat{y})$

# Optimal Classification Function

Expected loss is also called risk:

$$R(h) = \mathbb{E}_{XY}\big[L\big(Y, h(X)\big)\big]$$

$$h^* = \operatorname*{argmin}_{h} R(h)$$

$$= \operatorname*{argmin}_{h} \mathbb{E}_{XY}\big[L\big(Y, h(X)\big)\big]$$

{Whiteboard derivation}

# Optimal Classification Function

$$h^*(x) = \operatorname*{argmax}_{y \in \{0,1\}} P(Y = y \mid X = x)$$

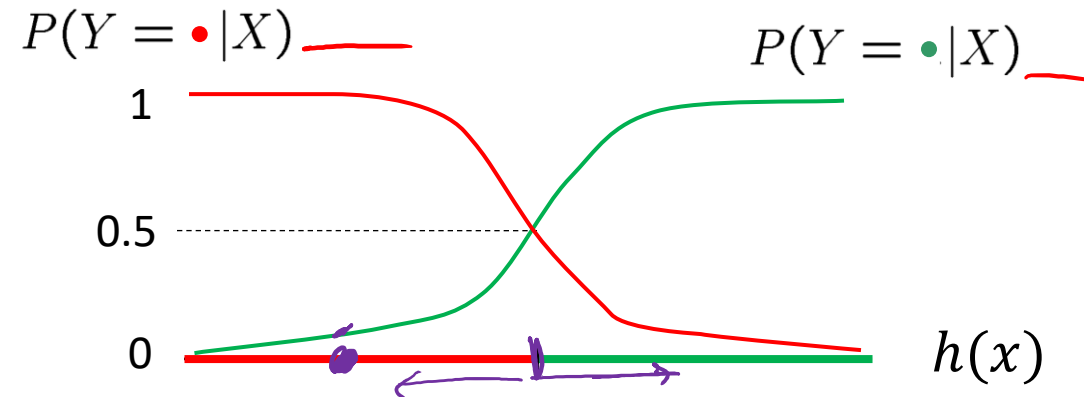$$h^* = \operatorname*{argmin}_{h} \mathbb{E}_{XY}\big[L\big(Y, h(X)\big)\big]$$

Start with arbitrary two-class loss $L(y, \hat{y})$

$$h^*(x) = \begin{cases} 1 & \text{if} \quad P(Y = 0 \mid x)L(0,1) + P(Y = 1 \mid x)L(1,1) \\ & \leq P(Y = 0 \mid x)L(0,0) + P(Y = 1 \mid x)L(1,0) \\ 0 & \text{otherwise} \end{cases}$$

Two-class, 0, 1 loss

$$h^*(x) = \begin{cases} 1 & \text{if} \quad P(Y = 0 \mid x) \\ & \leq P(Y = 1 \mid x) \\ 0 & \text{otherwise} \end{cases}$$

# Optimal Classification Function

$$h^*(x) = \underset{y \in \{0,1\}}{\mathrm{argmax}} \, P(Y = y \mid X = x)$$

$$h^* = \underset{h}{\mathrm{argmin}} \, \mathbb{E}_{XY}\big[L(Y, h(X))\big]$$



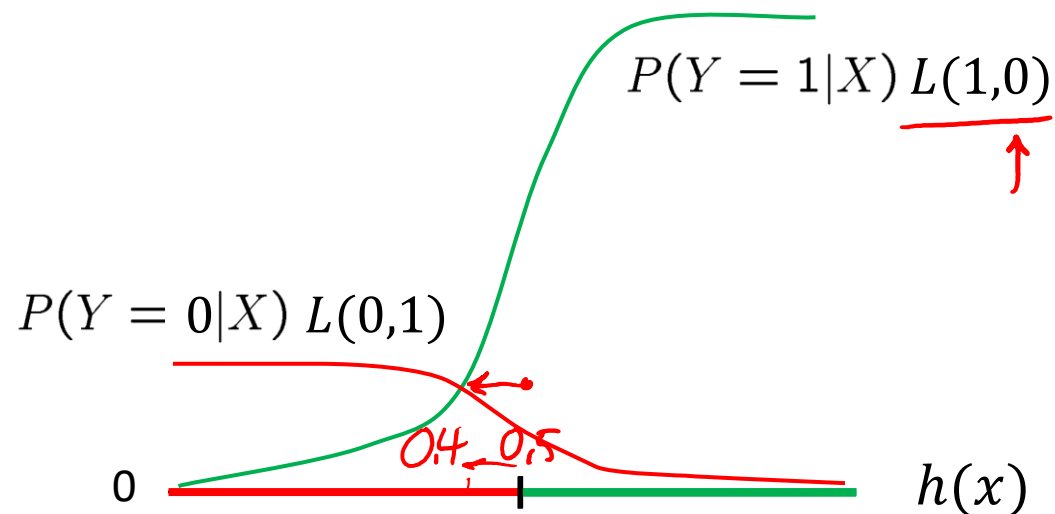**Start with arbitrary two-class loss $L(y, \hat{y})$**

$$h^*(x) = \begin{cases} 1 & \text{if} & P(Y = 0 \mid x)L(0,1) + P(Y = 1 \mid x)L(1,1) \\ & \leq & P(Y = 0 \mid x)L(0,0) + P(Y = 1 \mid x)L(1,0) \\ 0 & \text{otherwise} \end{cases}$$

**Two-class, weighted loss**

$$h^*(x) = \begin{cases} 1 & \text{if} & P(Y = 0 \mid x)L(0,1) \\ & \leq & P(Y = 1 \mid x)L(1,0) \\ 0 & \text{otherwise} \end{cases}$$

$P(Y = 1 \mid X)\,L(1,0)$

$P(Y = 0 \mid X)\,L(0,1)$

0.4  0.5

$0$       $h(x)$

# Optimal Classification Function

$$h^*(x) = \begin{cases} 1 & \text{if} & P(Y = 0 \mid x) \\ & \leq & P(Y = 1 \mid x) \\ 0 & \text{otherwise} \end{cases}$$

What is the risk of the optimal classifer?

$$R(h^*) = \mathbb{E}_{XY}\big[L\big(Y, h^*(X)\big)\big]$$

$$= \int f(x)\Big[ p(Y=0 \mid x) \underbrace{L(0, h(x))}_{h(x)} + p(Y=1 \mid x) \underbrace{L(1, h(x))}_{1-h(x)} \Big] dx$$

$$R(h^*) > 0$$



$P(Y = \textcolor{red}{\bullet} \mid X)$     $P(Y = \textcolor{green}{\bullet} \mid X)$

$R(h) - R(h^*)$

# Risk in Regression

Squared error loss $L(y, h(x)) = (h(x) - y)^2$

$$h^* = \underset{h}{\text{argmin}} \; \mathbb{E}_{XY}\big[L(Y, h(X))\big]$$

$$R(h) = E_{xy}\big[ L(Y, h(x)) \big]$$

$$= \iint (y - h(x))^2 f(dx, dy)$$

$$= E_X\left[ E_{Y|x}\big[ (Y - h(x))^2 \mid X \big] \right]$$

$$h^*(x) = \underset{\hat{y}}{\text{argmin}} \; E_{Y|x}\big[ (Y - \hat{y})^2 \mid X = x \big] \longrightarrow h^*(x) = E\big[ Y \mid X = x \big]$$

# Optimal Hypothesis Function

Goal: find a prediction function $h^*: \mathcal{X} \to \mathcal{Y}$ that minimizes the risk, the expected loss for randomly drawn test data $(X, Y)$

$$h^* = \operatorname*{argmin}_{h} R(h) = \operatorname*{argmin}_{h} \mathbb{E}_{XY}\big[L\big(Y, h(X)\big)\big]$$

# Learning from Training Data

But we want our hypothesis function to generalize well?

- How do we characterize and quantify this trade-off?

- {Back to the whiteboard}