$$D = \{(x_1, y_1) \cdots (x_n, y_n)\} \quad y_i \in \mathbb{R}$$

Each $(x, y)$ iid

$$P(x, y) = P(y|x) P(x)$$

y has distribution for same x

Expected label $\bar{y}(x) = E_{y|x}(y) = \int y \, P(y|x) \, dy$

Eq. Same house could be sold for $50,000 - 70,000\$.

Classifier being learned $h_D$

$$h_D = A(D) \qquad \leftarrow \text{ Algorithm}$$
say SVM, Perceptron

Expected test error given $h_D$

$$= E_{x,y \sim D} \left[ (h_D(x) - y)^2 \right] \qquad \text{Simplicity we pick square loss.}$$

$$= \int_x \int_y [h_D(x) - y]^2 P(x, y) \, dy \, dx$$

Now, $h_D$ is also a random variable set. As for different $D$, $h_D$ changes. Expected value of $h_D$

Expected classifier $= \bar{h} = E_{D \sim P^n} [A(D)] = \int_D h_D P(D) \, dD$

$\bar{h}$ — Average classifier on infinitely many datasets.

Expected Error of A

$$E_{(x,y) \sim P \atop D \sim P^n} \left[ (h_D(x) - y)^2 \right] \qquad \text{Take } D \text{ of } n \text{ datapts. from } P$$
Train to get $h_D(x)$, take a $(x, y)$ test point and get $(h_D(x) - y)^2$

$$= \int_D \int_x \int_y [h_D(x) - y]^2 P(x, y) P(D) \, dy \, dx \, dD$$

$$= E_{x,y \atop D} \left[ [\underbrace{h_D(x) - \bar{h}(x)}_{a} + \underbrace{\bar{h}(x) - y}_{b}]^2 \right] \qquad \text{Add an subtract } \bar{h}(x)$$

$$= \underbrace{E_{x,D} \left[ [h_D(x) - \bar{h}(x)]^2 \right]}_{a^2} + \underbrace{E_{x,y} \left[ [\bar{h}(x) - y]^2 \right]}_{b^2} + 2 \underbrace{E_{x,y \atop D} [(h_D(x) - \bar{h}(x))(\bar{h}(x) - y)]}_{\substack{2ab \approx 0 \\ (\text{see next page})}}$$

$$E_{x,y}\left[E_D\left[h_D(x) - \bar{h}(x)\right]\right]$$

$\downarrow$ because of linearity of Expected value ))

$2ab = 0$

$$E_D\left[h_D(x) - \bar{h}(x)\right] \qquad\qquad E_D\left[\bar{h}(x)\right] = \bar{h}(x)$$

$$\bar{h}(x) - \bar{h}(x) \qquad\qquad\qquad \uparrow \qquad\qquad \uparrow$$

$$\qquad\qquad\qquad\qquad\qquad\qquad D \text{ is not present here}$$

$$= 0$$

Thus, Expected Error of A remains as

$$E_{(x,y)\atop D}\left[(h_D(x) - y)^2\right] = E_{x,D}\left[\{h_D(x) - \bar{h}(x)\}^2\right]$$
$$+$$
$$E_{x,y}\left[(\bar{h}(x) - y)^2\right]$$

---

<u>Now,</u> $E\left[(\bar{h}(x) - y)^2\right] = E\limits_{x,y}\left[\left[\underbrace{(\bar{h}(x) - \bar{y}(x))}_{a} + \underbrace{(\bar{y}(x) - y)}_{b}\right]^2\right]$

add & subtract $\bar{y}(x)$

$$= E\left[(\bar{h}(x) - \bar{y}(x))^2\right] + E\left[(\bar{y}(x) - y)^2\right]$$

$$+ 2 E\limits_{x,y}\left[\underbrace{(\bar{h}(x) - \bar{y}(x))(\bar{y}(x) - y)}_{= 0}\right]$$

Since $E_{x,y}\left[(\bar{y}(x) - y)\right] = 0$ (see above explanation)
& below

$$E\limits_{x}\left[E\limits_{y|x}\left[\bar{y}(x) - y\right](\bar{h}(x) - \bar{y}(x))\right]$$ )))

To calculate
$E_{x,y}$ first we calculate $E_x$ and then $E_f$ y given $x$ "$E_{y|x}$"

So, $\bar{y}(x) - E\limits_{y|x}\left[y\right] = \bar{y}(x) - \bar{y}(x)$

$\uparrow$
$\bar{y}$ a constant as it's a mean.

We are ~~stop~~ left with
left

Expected Error of A

$$= \underset{x,d}{E}\left[\left[h_D(x) - \hat{h}(x)\right]^2\right] + \underset{x,d}{E}\left[\left[\bar{h}(x) - \bar{y}(x)\right]^2\right] + E\left[\left(\bar{y}(x) - y\right)^2\right]$$

<span style="margin-left:3em">Variance of classifier</span>    BIAS²    NOISE
of data

Variance of classifier
how much the classifier $h_D(x)$
vaires w.r.t to $\hat{h}(x)$ ie
average/expected classifier

Maybe add more
features

↓

If have unlimited data
and could get expected
classifier $\hat{h}(x)$ how much can $\hat{I}$
get to predict the expected label.
(not a particular $y$ but
average $y$)

If data is non linear and I use linear
classifier, no matter how much data
I have I will always have high
bias.