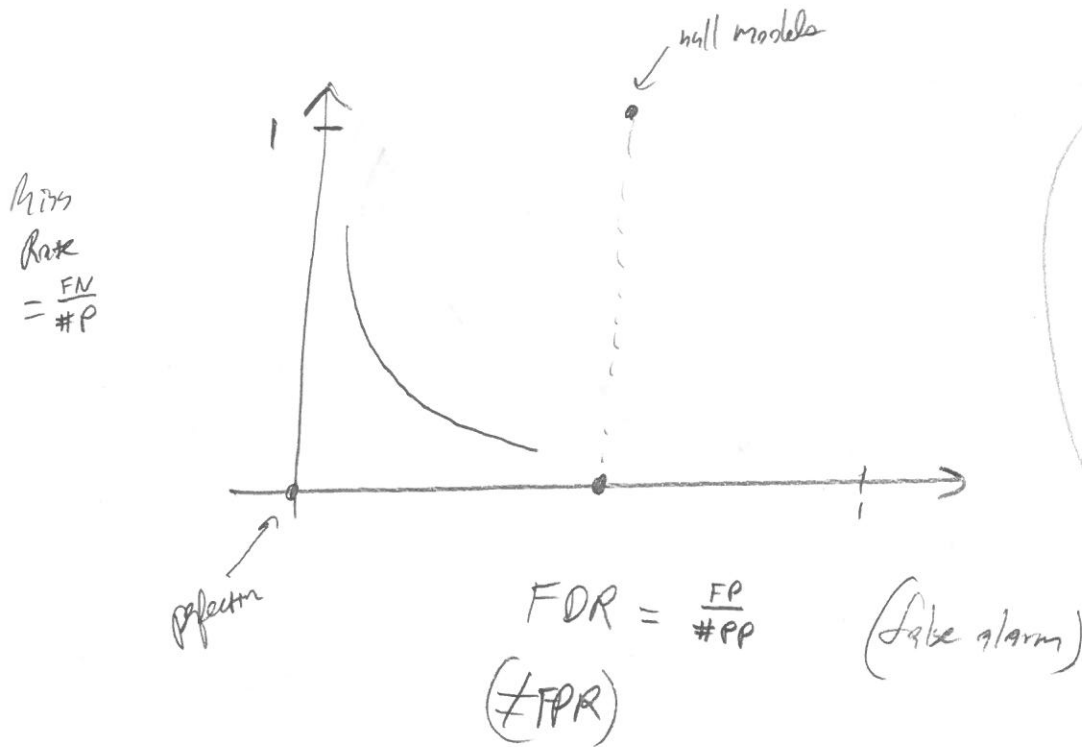


This is the most popular illustration of all model options, but not the only one. Also there's detection error tradeoff (DET, 1997)



Under the null model

$$FDR = \frac{p_{ch}(1-p_y)}{p_{ch}} = 1 - p_y \text{ (always)}$$

$$\text{Miss Rate} = \frac{(1-p_{ch})p_y}{p_y} = 1 - p_{ch}$$

if $p_{ch} = 0 \Rightarrow \text{all } \hat{y}_i = 1$ Never miss anything!

if $p_{ch} = 1 \Rightarrow \text{all } \hat{y}_i = 0$ miss everything FDR undefined

Review
 $\hat{y} = 1 \text{ if } p_{ch}$ param. which determines tradeoff of FPR & FN.
 Look at all choices of p_{ch} in a "performance curve"

ROC



AUC: evaluates the collection of classifier models from the prob. classifier \Rightarrow evaluate prob. classifier itself as a Brn

How to select a classifier - one pt on the ROC or DET plot! You need to define your costs/rewards

Bias-Variance Decomposition

Very theoretical...

Back to regression... recall

$$y = g + e = g + \underbrace{(f - g)}_{\substack{\text{error due to} \\ \text{missp. \& est} \\ \text{error}}} + \underbrace{\epsilon}_{\substack{\text{error due to} \\ \text{ignore \& est} \\ \text{error}}} \text{ aka "irreducible error"}$$

$$\Rightarrow \underbrace{y - g}_{\substack{\text{measure of} \\ \text{how good} \\ \text{the model } g \\ \text{is}}} = \underbrace{f - g} + \underbrace{\epsilon}$$

$$e^2 = (f - g + \epsilon)^2$$

What if I want the "mean squared error (MSE)". This would mean the expectation of the sq'd. residuals

In order to take an expectation, you need r.v.'s! Let's try this:

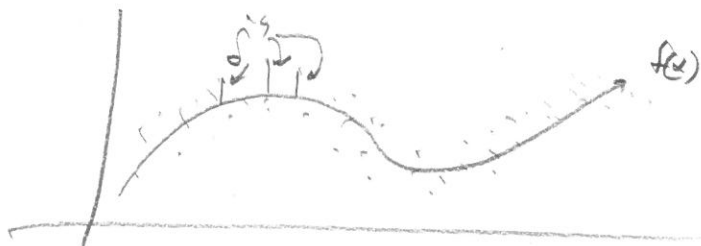
$$Y = f(\vec{x}) + \Delta \quad \begin{matrix} \nearrow \text{r.v. for response} & \nearrow \text{r.v. for irr. error} \end{matrix}$$

Realis really pure. How about...

$$Y | \vec{X} = \vec{x} = f(\vec{x}) + \Delta | \vec{X} = \vec{x}$$

Economists call this the "conditional expectation function" (CEF)

Let's assume $\textcircled{I} \quad \underbrace{E[Y | \vec{X} = \vec{x}]}_{\substack{\uparrow \\ \text{avg over irr. err}}} = f(\vec{x}) \Rightarrow E[\Delta | \vec{X} = \vec{x}] = 0$



so... for any \vec{x} , if many realizations of $Y \Rightarrow$ the avg of those Y 's is close to $f(\vec{x})$.

Variance does not depend on \vec{x} (homoskedastic assumption)

(F)

$$\textcircled{I} \text{Var}[\Delta | \vec{X} = \vec{x}] = \text{Var}[\Delta] = \sigma^2 \Rightarrow E[\Delta^2] = \sigma^2$$

Back to MSE, we care about this for a new obs. \vec{x}^* . This is also called "generalization error."

$$\text{MSE}(\vec{x}^*) := E[(Y^* - g(\vec{x}^*))^2 | \vec{X} = \vec{x}^*] \geq \sigma^2$$

Can't be lower than irreducible error

What if we know f ?

$$\text{MSE}(\vec{x}^*) := E[(Y^* - f(\vec{x}^*))^2 | \vec{X} = \vec{x}^*] = E[\Delta^{*2} | \vec{X} = \vec{x}^*] = E[\Delta^2] = \sigma^2$$

What is this expectation over? Δ^* . An integral over \mathbb{R} with $P(\Delta)$ inside

If f unknown...

$$\begin{aligned} \text{MSE}(\vec{x}^*) &= E[Y^{*2} - 2Y^*g(\vec{x}^*) + g(\vec{x}^*)^2 | \vec{X} = \vec{x}^*] \\ &= E[Y^{*2} | \vec{X} = \vec{x}^*] - 2 E[Y^*g(\vec{x}^*) | \vec{X} = \vec{x}^*] + E[g(\vec{x}^*)^2 | \vec{X} = \vec{x}^*] \\ &= E[(f(\vec{x}) + \Delta)^2 | \vec{X} = \vec{x}^*] \quad E[Y^*]g(\vec{x}^*) \quad g(\vec{x}^*)^2 \\ &= E[(f(\vec{x}) + \Delta)^2] \quad \parallel \quad f(\vec{x}^*)g(\vec{x}^*) \\ &= E[f(\vec{x})^2 + 2f(\vec{x})\Delta + \Delta^2] \\ &= f(\vec{x})^2 + \sigma^2 - 2f(\vec{x}^*)g(\vec{x}^*) + g(\vec{x}^*)^2 \\ &= (f(\vec{x}^*) - g(\vec{x}^*))^2 + \sigma^2 \end{aligned}$$

Expected sqd error is additive

$$e = (f - g) + f$$

Not so interesting... How about we take expectation over all $\Delta_1, \Delta_2, \dots, \Delta_n$ and Δ^* ? Expectation over all \mathcal{D} !

$$MSE(\vec{x}^*) = E_{\Delta_1, \dots, \Delta_n, \Delta^*} [Y^2 | \vec{X} = \vec{x}^*] - 2 E_{\Delta_1, \dots, \Delta_n, \Delta^*} [Y^* g(\vec{x}^*) | \vec{X} = \vec{x}^*] + E_{\Delta_1, \dots, \Delta_n} [g(\vec{x}^*)^2 | \vec{X} = \vec{x}^*]$$

only varies with Δ^*

$$A(D, \mathcal{H}) = R(X, \vec{y}, \mathcal{H})$$

which is a function of \mathcal{D} 's!
but not a function of \mathcal{H}

$$= E_{\Delta^*} [Y^2 | \vec{X} = \vec{x}^*] - 2 E_{\Delta_1, \dots, \Delta_n, \Delta^*} [Y^* g(\vec{x}^*) | \vec{X} = \vec{x}^*] + E_{\Delta_1, \dots, \Delta_n} [g(\vec{x}^*)^2]$$

↑ ↑
independent

$$= \sigma^2 + f(\vec{x}^*)^2 - 2 \underbrace{E_{\Delta^*} [Y^*]}_{f(\vec{x}^*)} E_{\Delta_1, \dots, \Delta_n} [g(\vec{x}^*)] + E_{\Delta_1, \dots, \Delta_n} [g(\vec{x}^*)^2]$$

$$= \sigma^2 + Var[g(\vec{x}^*)] + f(\vec{x}^*)^2 - 2 f(\vec{x}^*) E[g(\vec{x}^*)] + E[g(\vec{x}^*)]^2$$

$$= \sigma^2 + Var[g(\vec{x}^*)] + (E[g(\vec{x}^*)] - f(\vec{x}^*))^2$$

$$= \sigma^2 + Var[g(\vec{x}^*)] + \underbrace{E[g(\vec{x}^*) - f(\vec{x}^*)]^2}_{Bias[g(\vec{x}^*)]^2}$$

How much g
varies around its mean
function $E[g(\vec{x})]$
(stability in learning)

Bias-Var tradeoff

How far away the avg g is
from f .

One more small pt. This was all for one data pt \vec{x}^* . We can average over all $\vec{x}_1, \dots, \vec{x}_n, \vec{x}^* \in \mathcal{X}$. Assume distr. $P(\vec{X})$.

indep of \mathcal{X}

$$MSE = E_{\mathcal{X}} [\sigma^2 + Var[g(\vec{x}^*)] + Bias[g(\vec{x}^*)]^2]$$

$$= \sigma^2 + E_{\mathcal{X}} [Var[g(\vec{x}^*)]] + E_{\mathcal{X}} [Bias[g(\vec{x}^*)]^2]$$

Since this one \mathcal{D} and
one g ... we don't
get to see this
except during
simulation

gen. error.

expected
variance

expected
Bias²

deno