# CPH 100A: Machine Learning Foundations I

**Instructor:** Adam Yala, PhD ([yala@berkeley.edu](mailto:yala@berkeley.edu))

**Computational**
PRECISION HEALTH

Berkeley | UCSF

# Problem Motivation: Early Detection is critical



**5-year survival** / **Stage** / **Lung Cancer**

Legend: Localized, Regional, Distant



Size of cancer / Time

- Fast
- Slow
- Very slow
- Nonprogressive

Abnormal cell — Screening detects cancer — Death from other causes

Size at which cancer causes death

Size at which cancer causes symptoms

This is when overdiagnosis occurs

# RCTs reduce lung cancer mortality

Reduced Lung-Cancer Mortality with Low-Dose Computed Tomographic Screening

The National Lung Screening Trial Research Team
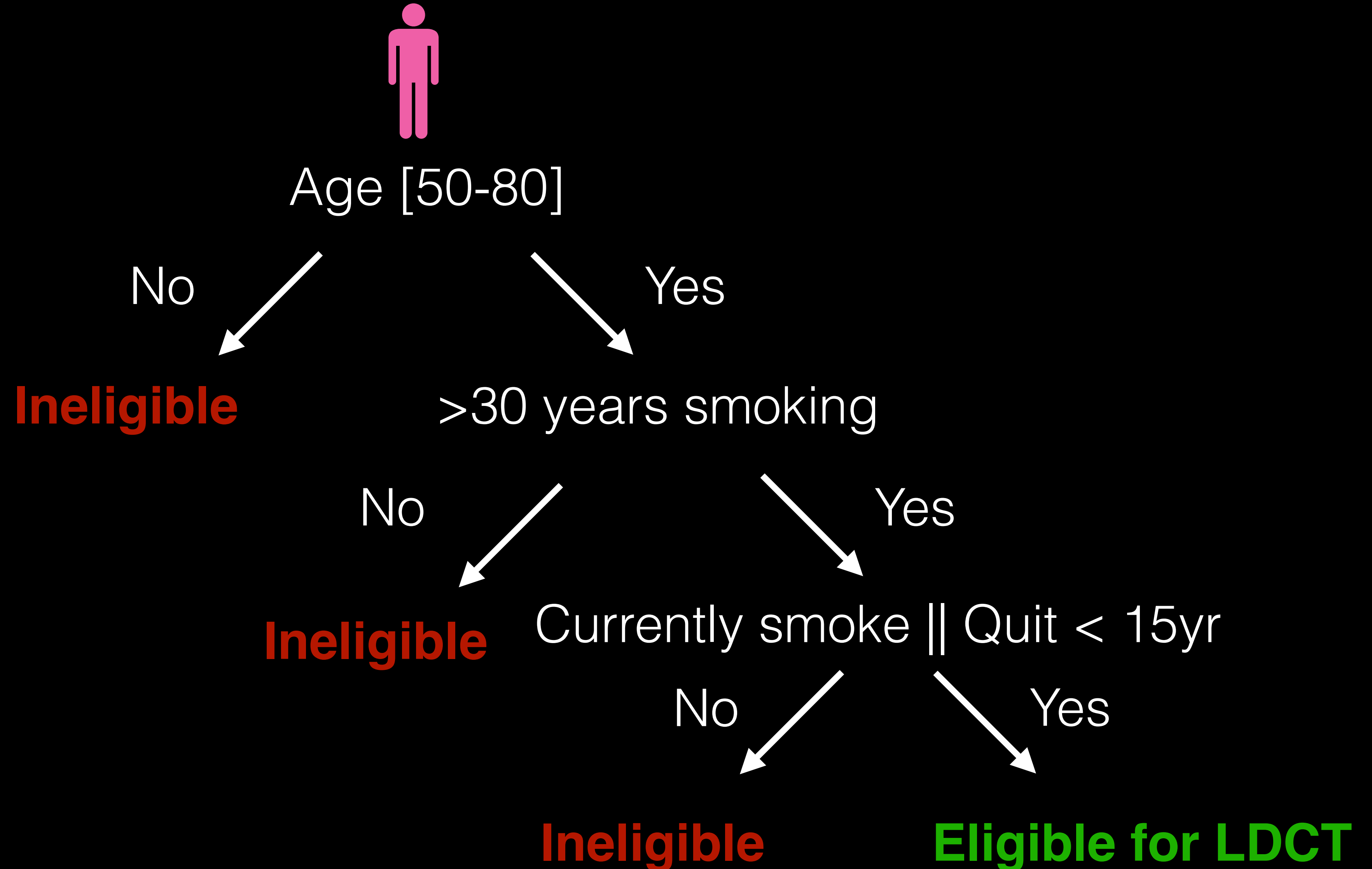
## NLST reduces lung cancer mortality by 20%

Reduced Lung-Cancer Mortality with Volume CT Screening in a Randomized Trial

Harry J. de Koning, M.D., Ph.D., Carlijn M. van der Aalst, Ph.D., Pim A. de Jong, M.D., Ph.D., Ernst T. Scholten, M.D., Ph.D., Kristiaan Nackaerts, M.D., Ph.D., Marjolein A. Heuvelmans, M.D., Ph.D., Jan-Willem J. Lammers, M.D., Ph.D., Carla Weenink, M.D., Uraujh Yousaf-Khan, M.D., Ph.D., Nanda Horeweg, M.D., Ph.D., Susan van 't Westeinde M.D., Ph.D., Mathias Prokop, M.D., Ph.D., et al.

## NELSON reduces lung cancer mortality by 24%

# NLST screening criteria



Age [50-80]

No → **Ineligible**

Yes → >30 years smoking

No → **Ineligible**

Yes → Currently smoke || Quit < 15yr

No → **Ineligible**
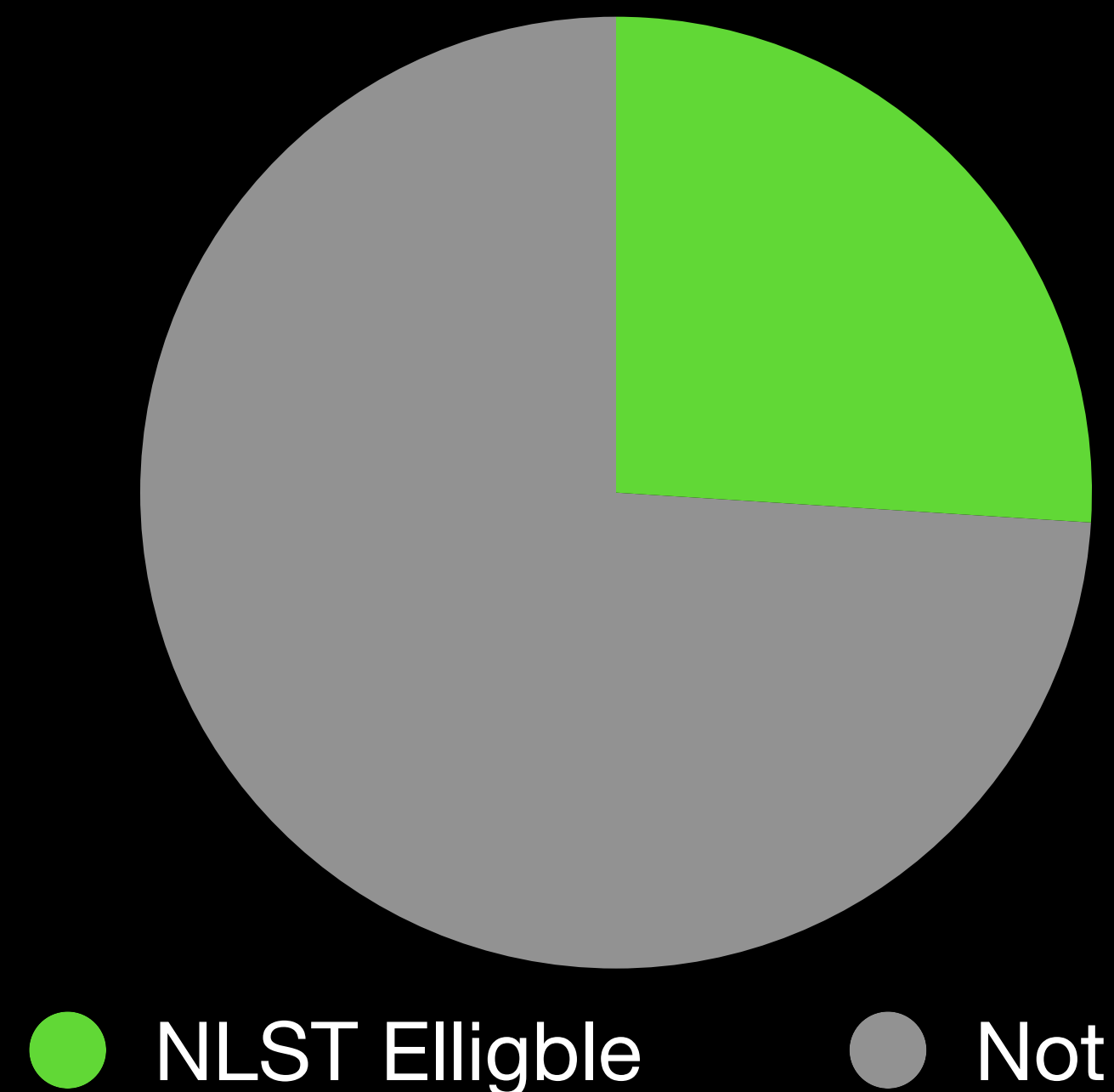
Yes → **Eligible for LDCT**

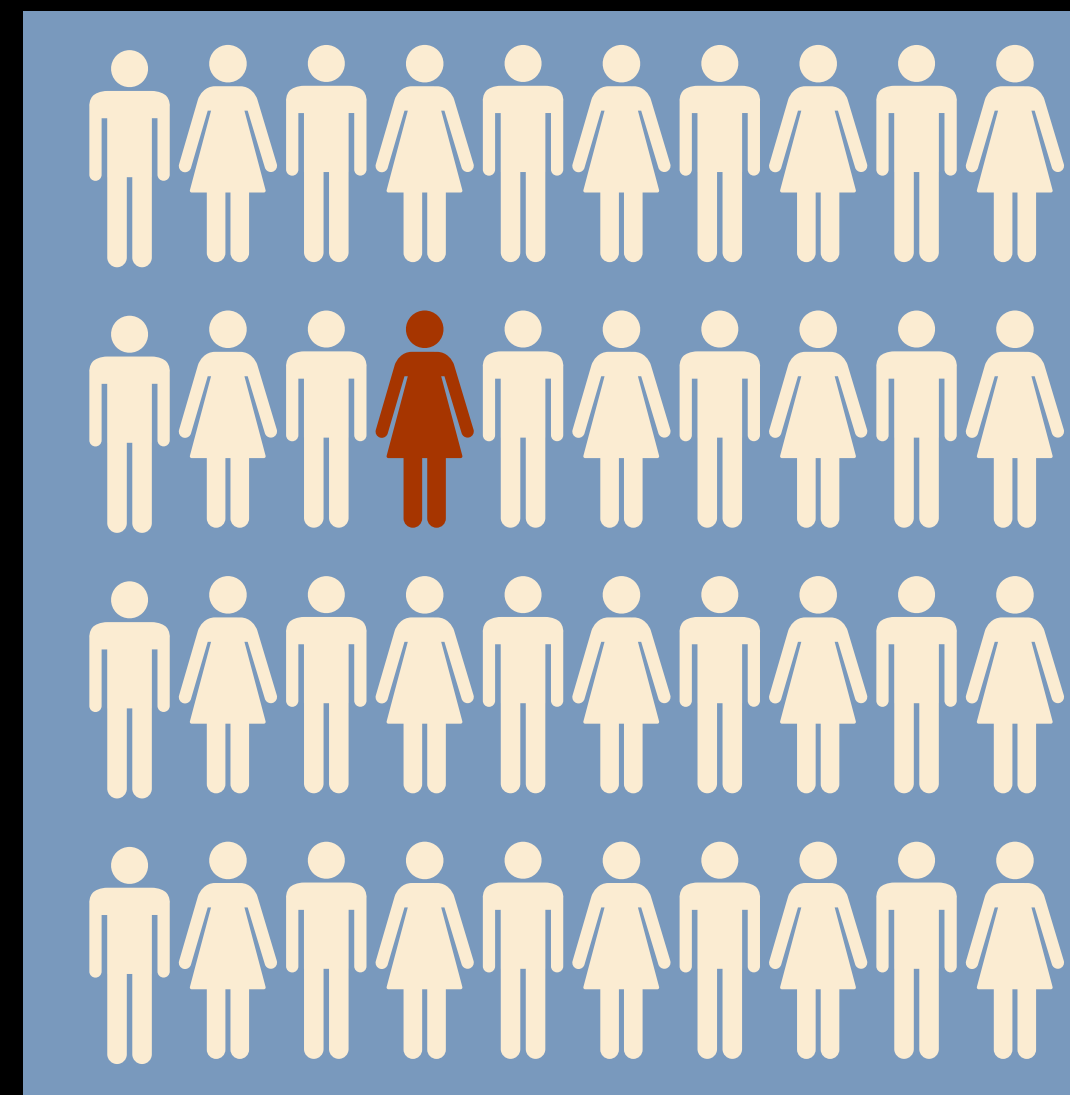# Efficacy of a screening program

Fundamental challenge is **cost-effectiveness**

How much benefit does it achieve?

How much harm does the program do?



NLST Elligble    Not

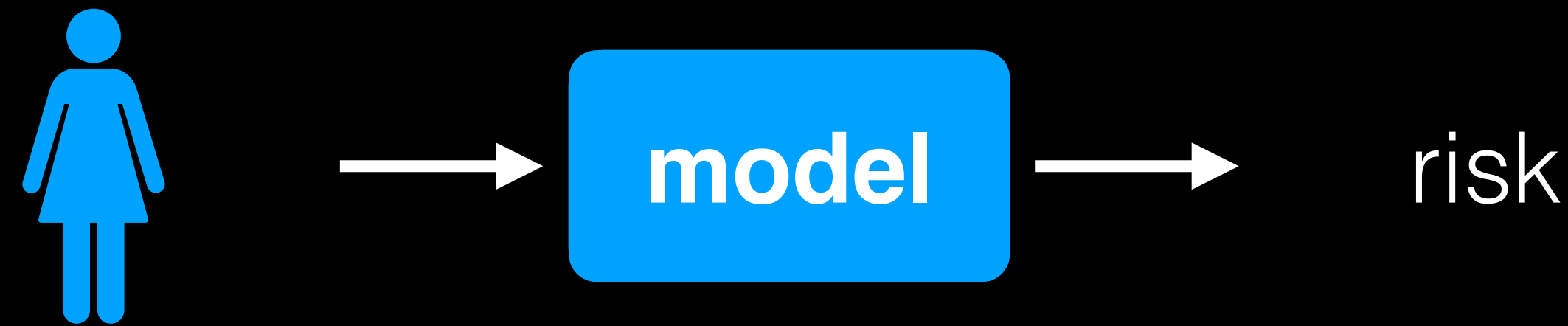PMID: 23060474

1000 screens

↓

240 positives

↓

6 cancers

# Can we do better?



Predict probability of cancer *(proxy for prob screening benefit)*

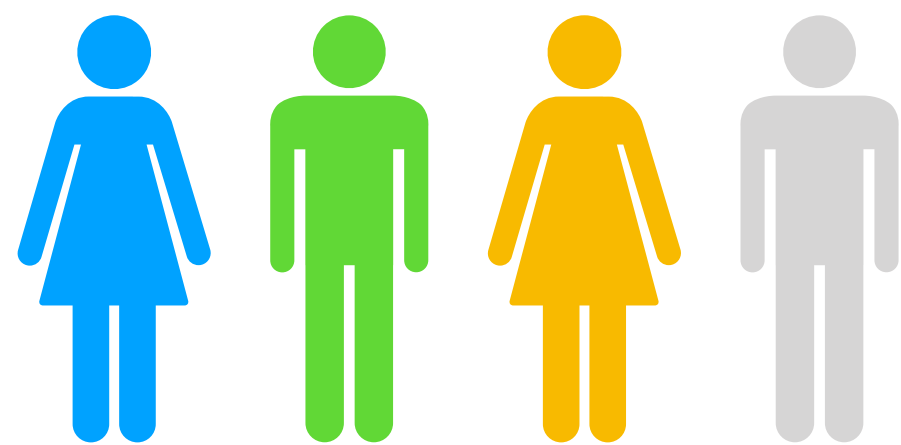Identify population with **higher specificity** and **higher sensitivity**

**Key Question for today:**

How do these models work?

How should we be evaluating them?

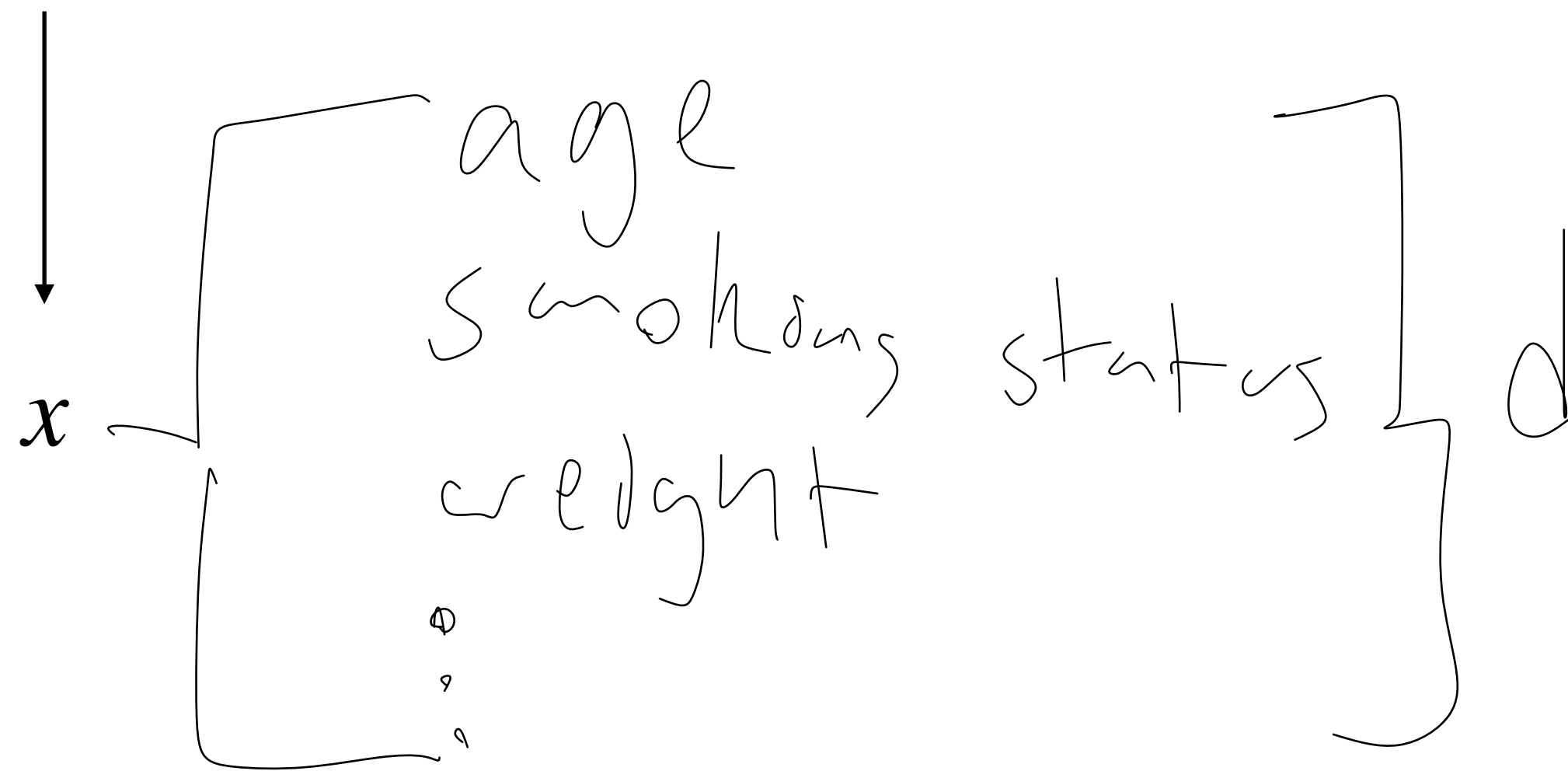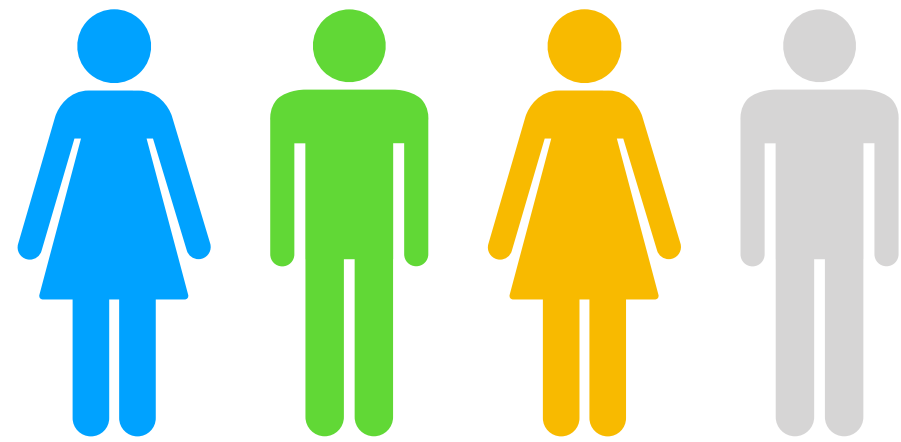# What we know

$n$ historical  patients



$x$                          $x \in \mathscr{R}^D$

$y$                          $y \in \{0,1\}$

# What we know

$n$ historical  patients



$$x \begin{bmatrix} \text{age} \\ \text{smoking status} \\ \text{creight} \\ \vdots \end{bmatrix} d$$

$$x \in \mathscr{R}^D$$

$$y \rightarrow cancer$$

$$y \in \{0,1\}$$

# What we want

$x \longrightarrow$ **model:** $h$ $\longrightarrow y$

# What we want

$x \longrightarrow$ **model:** $h$ $\longrightarrow y$

Some func $h$ $x \rightarrow h \rightarrow y$

and $h$ to be " good "

What is $h$ ?

$h : \mathbb{R}^D \rightarrow \mathbb{R}$

# Todays Hypothesis class: linear models
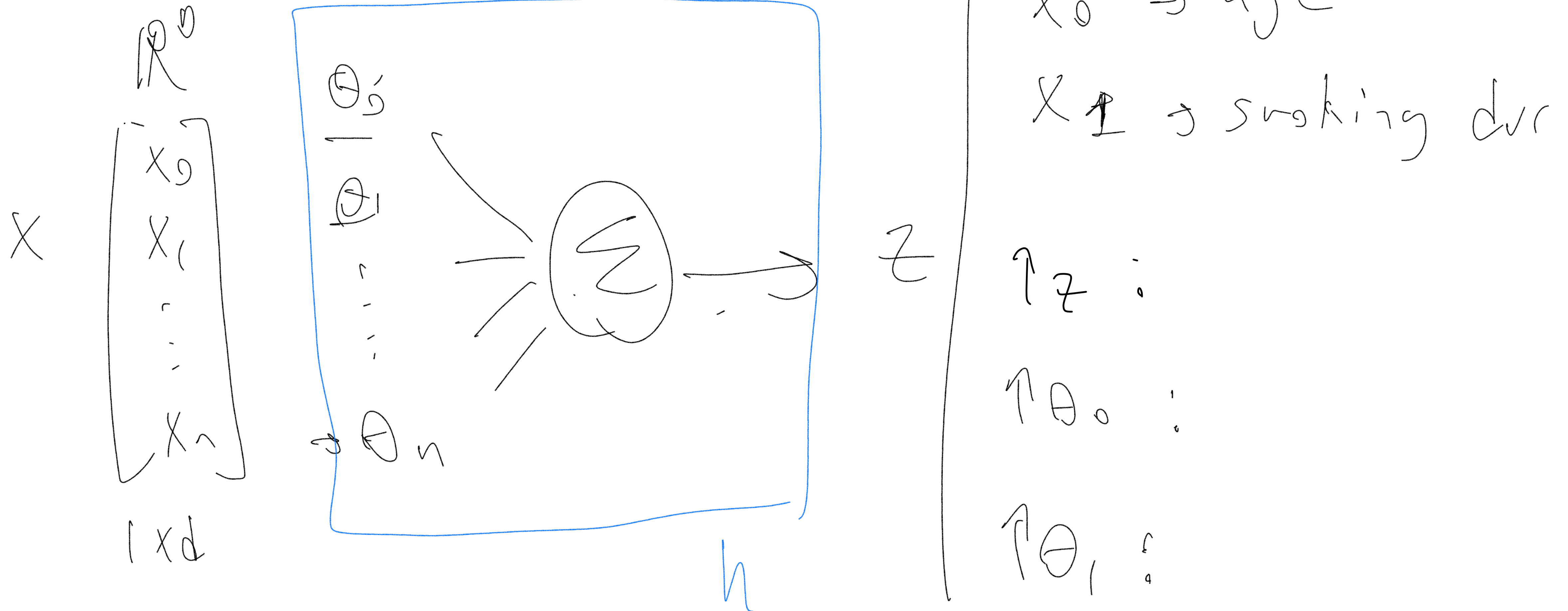
# Todays Hypothesis class: linear models

$\mathbb{R}^D$

$\mathbb{R}^p$

$X$

$\begin{bmatrix} X_0 \\ X_1 \\ \vdots \\ X_n \end{bmatrix}$

$1 \times d$

$\Theta_0$

$\Theta_1$

$\cdots$

$\to \Theta_n$

$\Sigma$

$\to z$

$1 \times d$

$h$

$$z = \sum_i^N \Theta_i X_i + b$$

$$z = \Theta X^T$$

simplified notation

# Interpreting Linear Models

$\mathbb{R}^D$

$X$

$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$

$1 \times d$



$\theta_0$

$\theta_1$

$\sum$

$\to \theta_n$

$h$

$z$

$x_0 \to age$

$x_1 \to smoking \; dur$

$\uparrow z$ :

$\uparrow \theta_0$ :

$\uparrow \theta_1$ :

# Interpreting Linear Models



$X$

$$\mathbb{R}^D$$

$$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$1 \times d$

$\theta_0$

$\theta_1$

$\sum$

$\theta_n$

$h$

$Z$

$x_0 \to$ age

$x_1 \to$ smoking dur

$\uparrow z$ : higher cancer risk

$\uparrow \theta_0$ : age pos assoc

$\uparrow \theta_1$ : smoking pos assoc

# Geometric View

# Capturing Uncertainty

# Capturing Uncertainty

$$\mathcal{H} : \theta X^T$$
$$R^D \to R$$

but we want

$$R^P \to [0,1]$$

we want:

$p$



$z$

$\uparrow$ scores $\Rightarrow p(x) = 1$

$\downarrow$ score $\Rightarrow p(x) = 0$

# Loglinear Models

# Loglinear Models

$$P = \sigma(z)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$z$

$$h(x) = \sigma(\theta x^T)$$

$P$   $P$   $P$

$x_0$   $x_0$   $x_0$

1D case

# Geometric View

# Empirical Risk Minimization

How do fnd a "good" $h$?

# Empirical Risk Minimization

How do find a "good" $h$?



vs

How good are the probs?

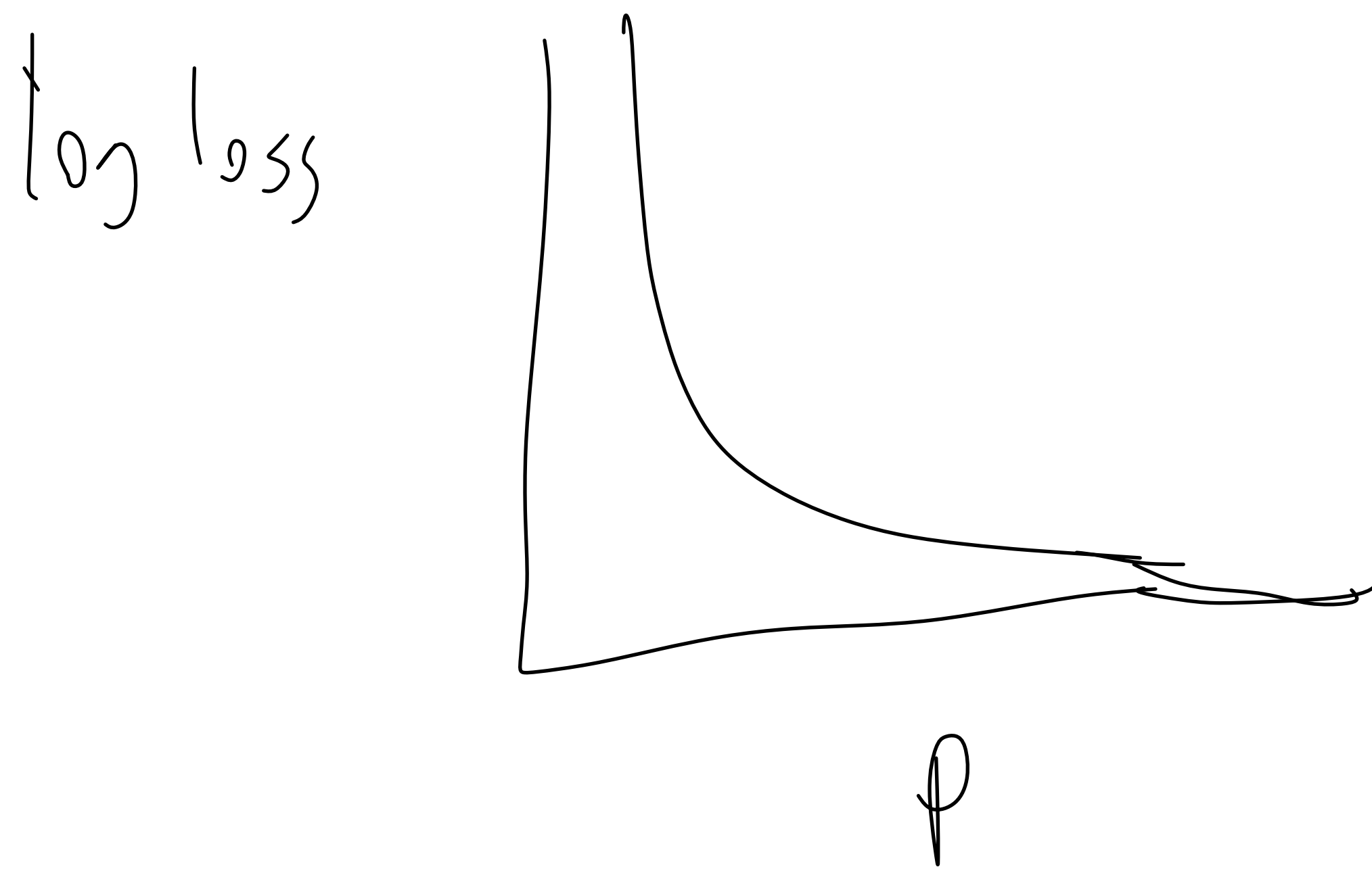$$R(\Theta) = \frac{1}{n} \sum_{i}^{N} L(y, p) = \frac{1}{n} \sum_{i}^{N} L(y, \sigma(\Theta x^T))$$

# Loss Function: Cross Entropy

# Loss Function: Cross Entropy

Likelihood of observed data

$$L(y, p) = -(y \log p + (1-y) \log(1-p))$$

$y \approx 1$

log loss
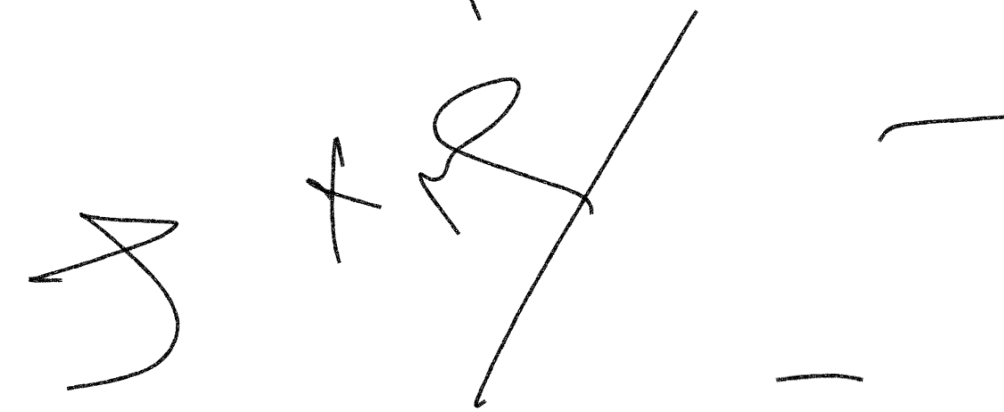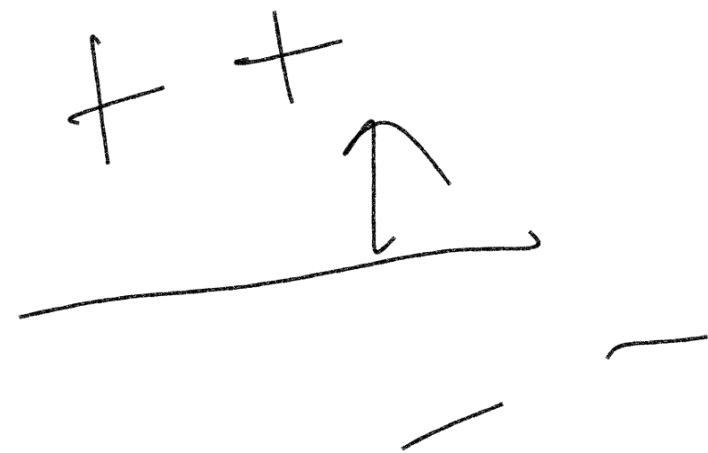


$p$

penalize very
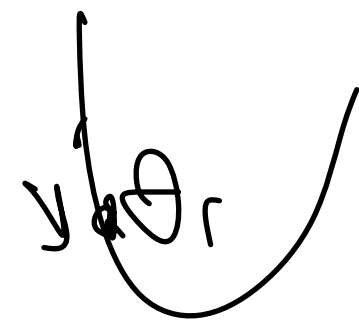wrong probs.

Maximum Likelihood Estimation!

# Optimization

How do we find a good $h$ ?

# Optimization

How do we find a good $h$?

• gradient descent!



$\theta_0$

$\theta_1$

$\theta_2$

# How? Gradient Descent

# How? Gradient Descent

$$\Theta_{new} = \Theta_{old} - \eta \frac{\partial R(\Theta_{old})}{\partial \Theta_{old}}$$

$\eta$ ← learning rate

$$\frac{\partial R(\Theta)}{\partial \Theta} = \frac{\partial}{\partial \Theta}\left(\frac{1}{n}\sum_i^N L(y,p)\right) = \frac{1}{n}\sum_i^N \frac{\partial L(y_i,p_i)}{\partial \Theta}$$

$$= \frac{1}{n}\sum_i^N (p_i - y_i)\, x_i$$

$$\frac{\partial L}{\partial \theta}(y, p) = \frac{\partial}{\partial \theta}\left(\right)$$

$$\frac{\partial L}{\partial \theta}(y, p) = \frac{\partial}{\partial \theta}\left( y \log \underbrace{\sigma(\theta x^T)}_{p} + (1-y)\log(1 - \sigma(\theta x^T)) \right)$$

$$= \left( \frac{y(1-p)\,p\,x^T}{p} + \frac{(1-y)\,p(1-p)(x)}{1-p} \right.$$

$$\left. = \left( y x - pyx - px + pyx \right) \right.$$

$$= (y-p)x \qquad\qquad = (p-y)x$$

# Putting it all together

$$\sigma\left(\theta_i x_i^T\right)$$

$$\downarrow$$

$$\theta_{new} = \theta_{old} - \eta \frac{\sum_{i}^{N}\left(P_i - y_i\right) x_i}{N}$$

# Putting it all together

$$\Theta_i x_i^T$$

$$\Theta_{new} = \Theta_{old} - \eta \frac{\sum_i^N (P_i - y_i) x_i}{N}$$

$\Theta_0$ init rand

while not converged:

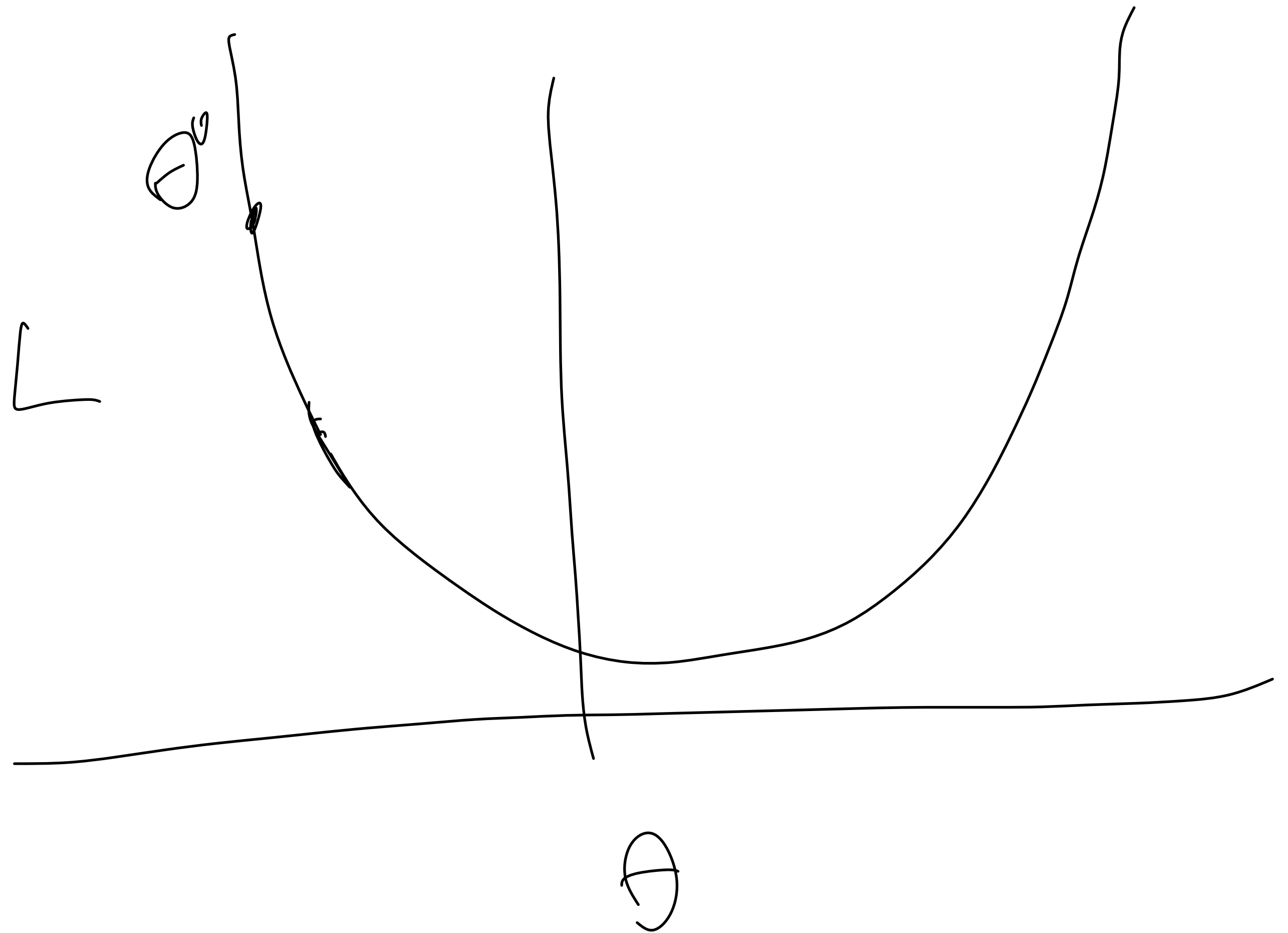$$\Theta_{t+1} = \Theta_t - \eta \frac{\partial L}{\partial \Theta_t}$$

return $\Theta_{end}$

# What if N is too large?

# What if N is too large?

## Stochastic Gradient Descent

Estimate empirical risk $R$
with $B$ random samples

$$\theta_i x_i^T$$

$$\Downarrow$$

$$\theta_{new} = \theta_{old} - \eta \frac{\sum_i^B (P_i - y_i) x_i}{B}$$

$\theta_0$ init rand

while not converged:

$$\theta_{t+1} = \theta_t - \eta \frac{\partial L}{\partial \theta_t}$$

return $\theta_{end}$

# Choosing your learning rate

What if LR is too small?

What if LR is too large?

How can we tell?

# Choosing your learning rate

What if LR is too small?

Slow optim

What if LR is too large?

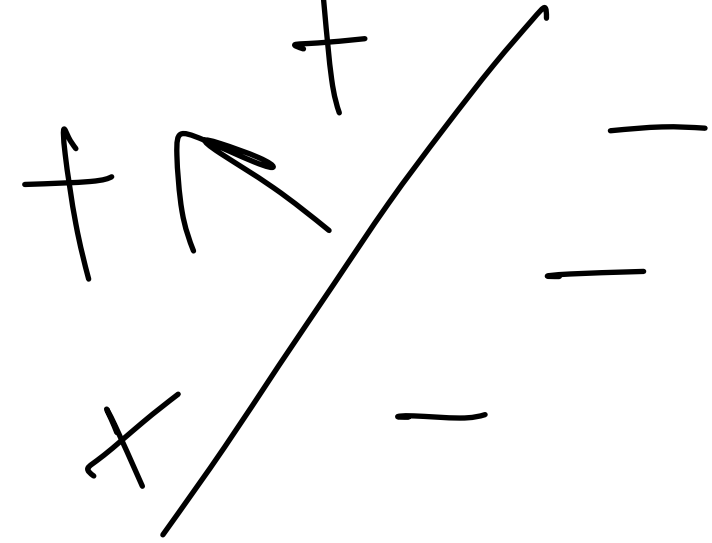Diverge

How can we tell?

Training Curve
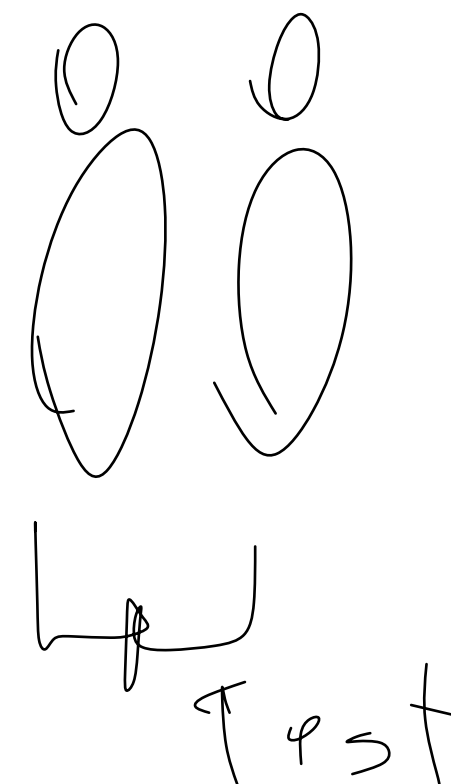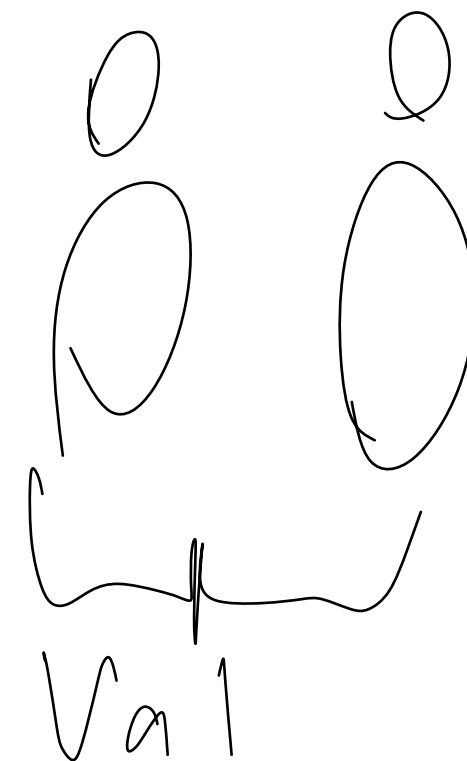
# Recap

# Recap

Now we have a model $h$ w param $\theta$
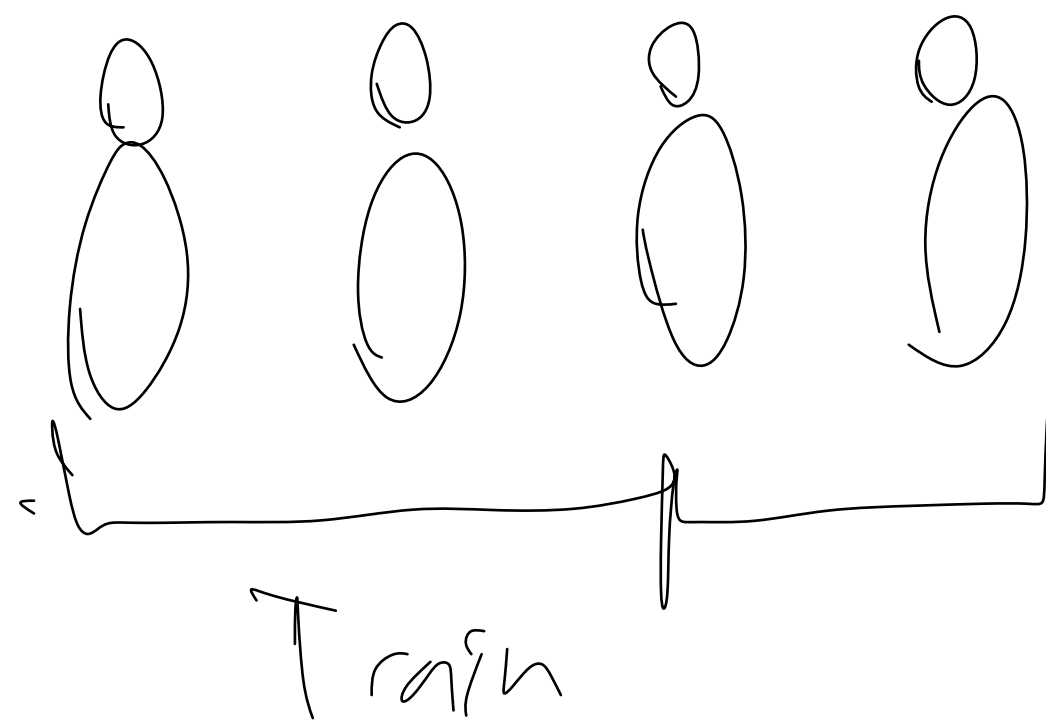
$$\theta = \arg\min_{\theta} R(\theta) = \sum_i^n \frac{L(\rho_i, \hat{y}_i)}{n}$$

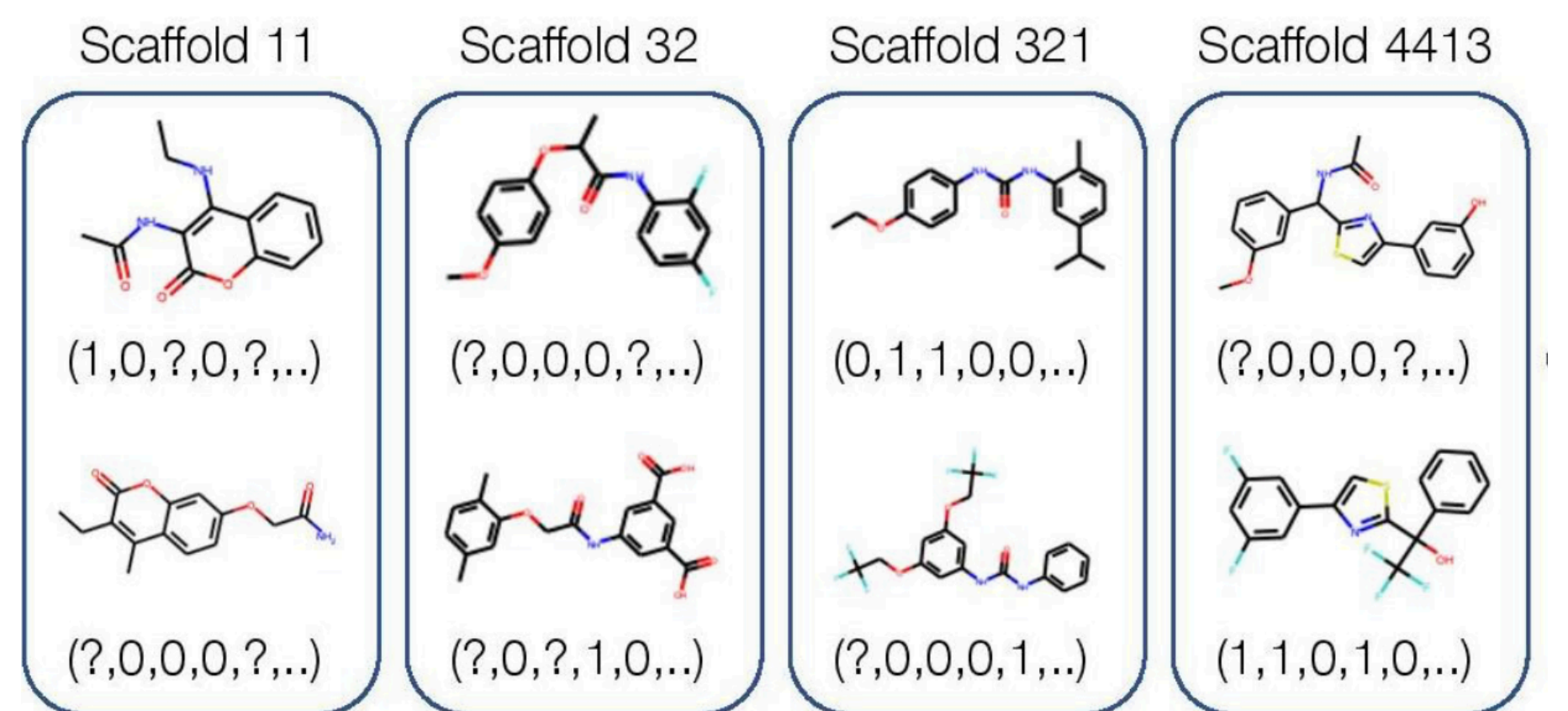is this enough?    No

We want $h$ to do well on NEW patients ...
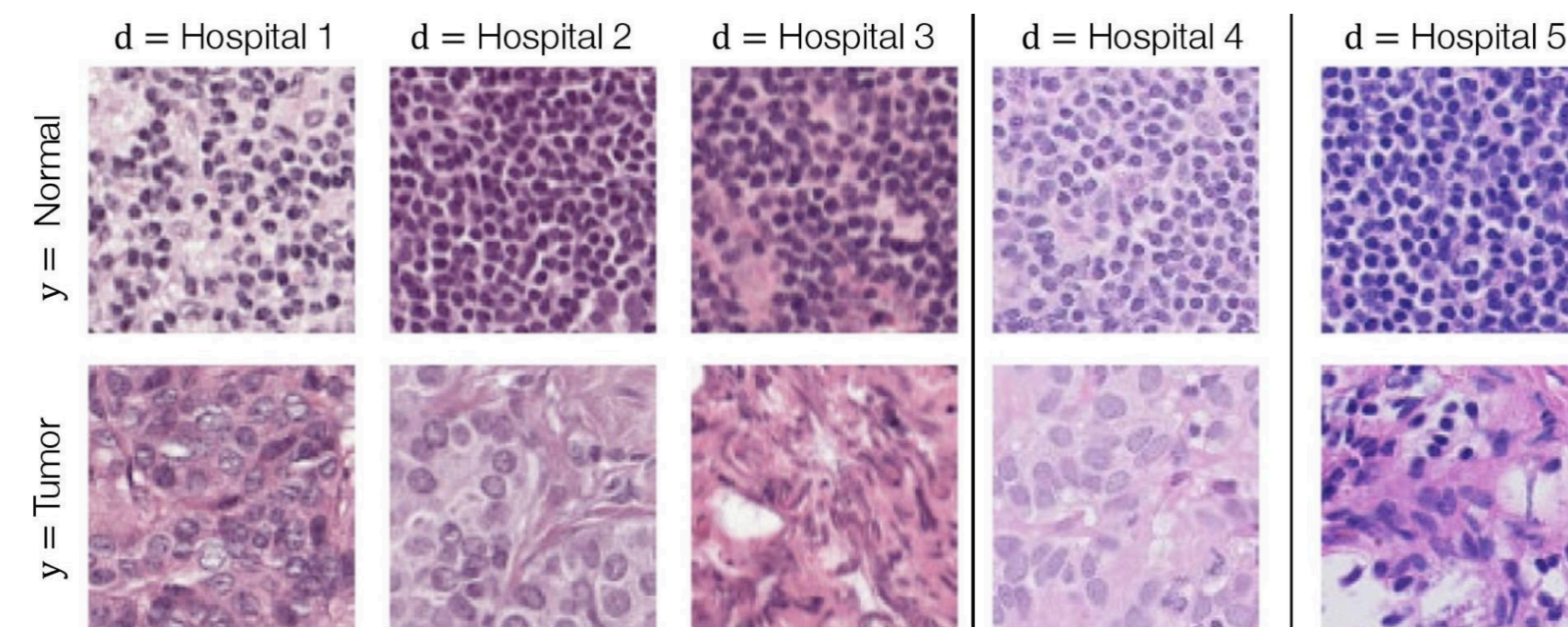


Train          Val          Test

# The importance of data splitting
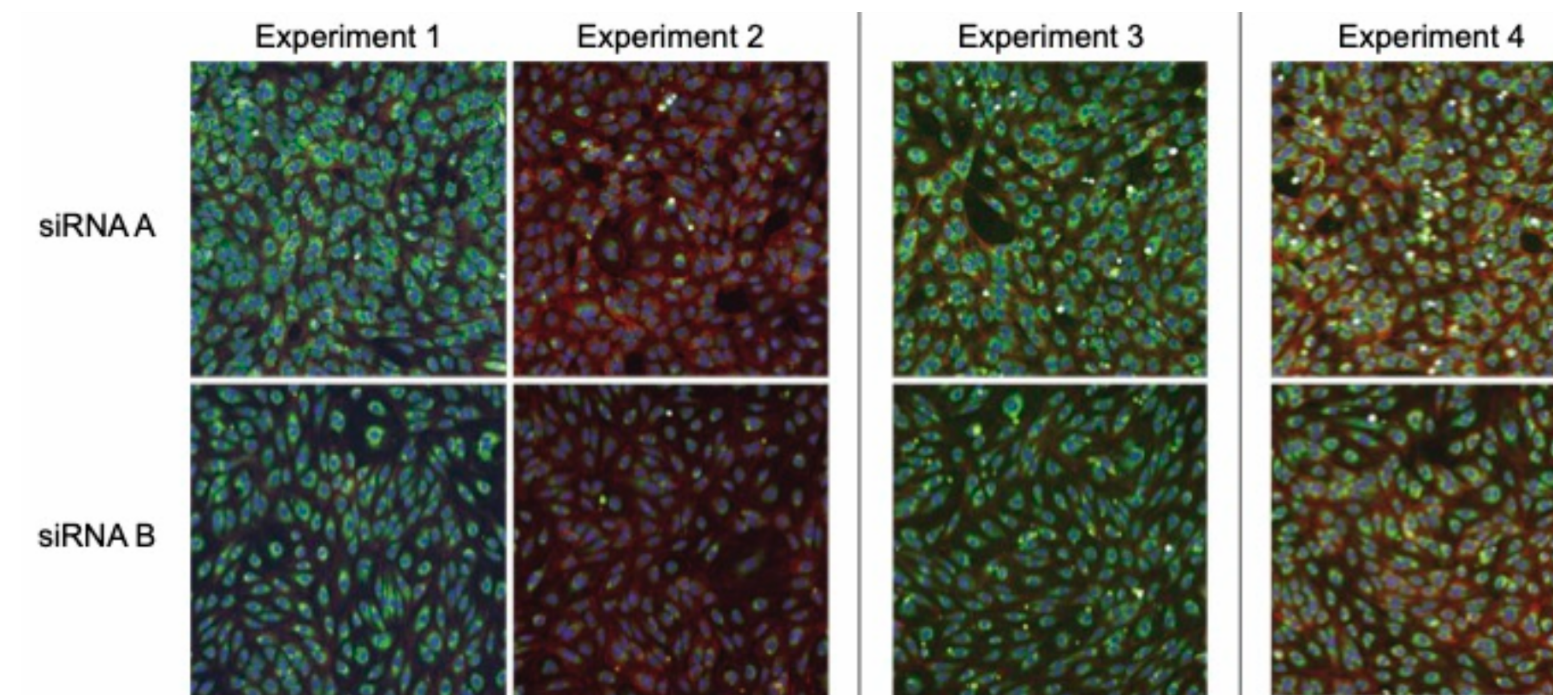
## Scaffold split in property prediction



Yang, Kevin, et al. "Analyzing learned molecular representations for property prediction." *Journal of chemical information and modeling* 59.8 (2019): 3370-3388.
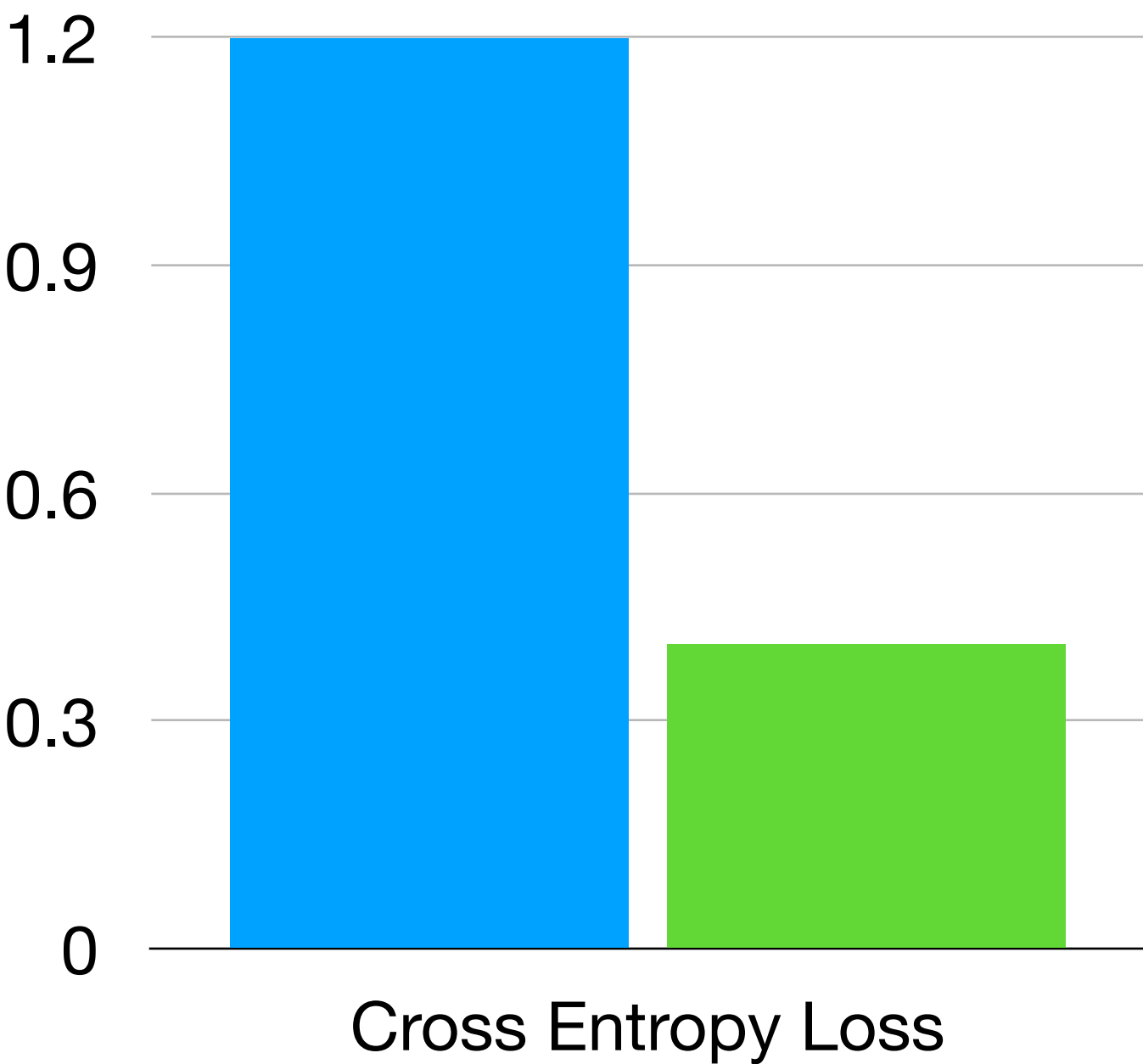
## Hospital source in pathology



Bandi, Peter, et al. "From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge." *IEEE transactions on medical imaging* 38.2 (2018): 550-560.

## Batch effects in high-throughput screening



Taylor, J., et al. "RxRx1: An Image Set for Cellular Morphological Variation Across Many Experimental Batches." *The 7th International Conference on Learning Representations*. 2019.

# Model Evaluation



1.2
0.9
0.6
0.3
0
Cross Entropy Loss

Modeling objective

Achievable performance

$1 : h(x) \geq p$

$0 : h(x) < p$

$TPR = \dfrac{TP}{\#+}$

$FPR = \dfrac{FP}{\#-}$

$AUC : P(P_i > P_j \mid y_i = 1, y_j = 0)$

**Computational**
PRECISION HEALTH

# Model Evaluation


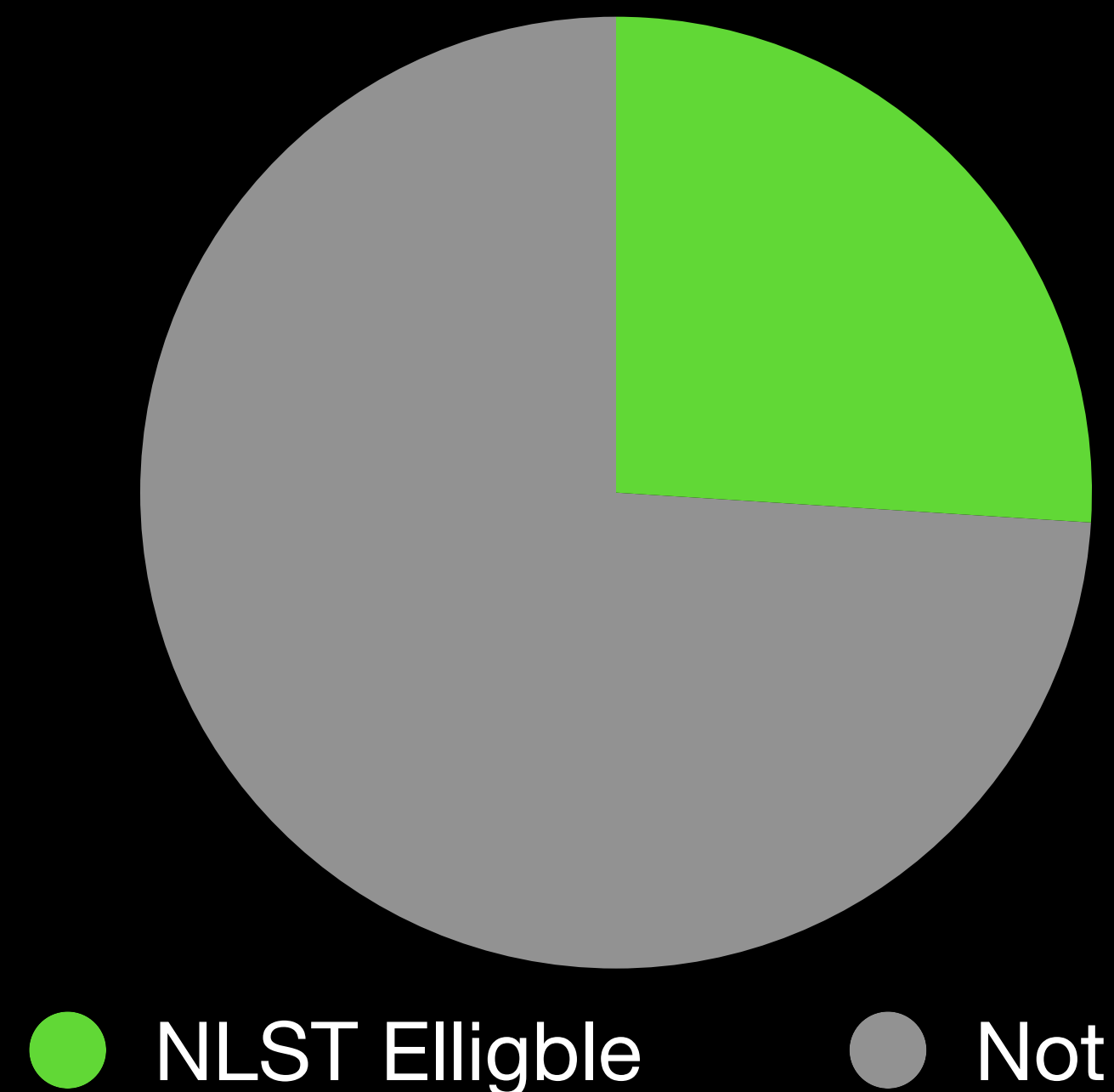
Modeling objective

Achievable performance

**Simulated clinical utility**
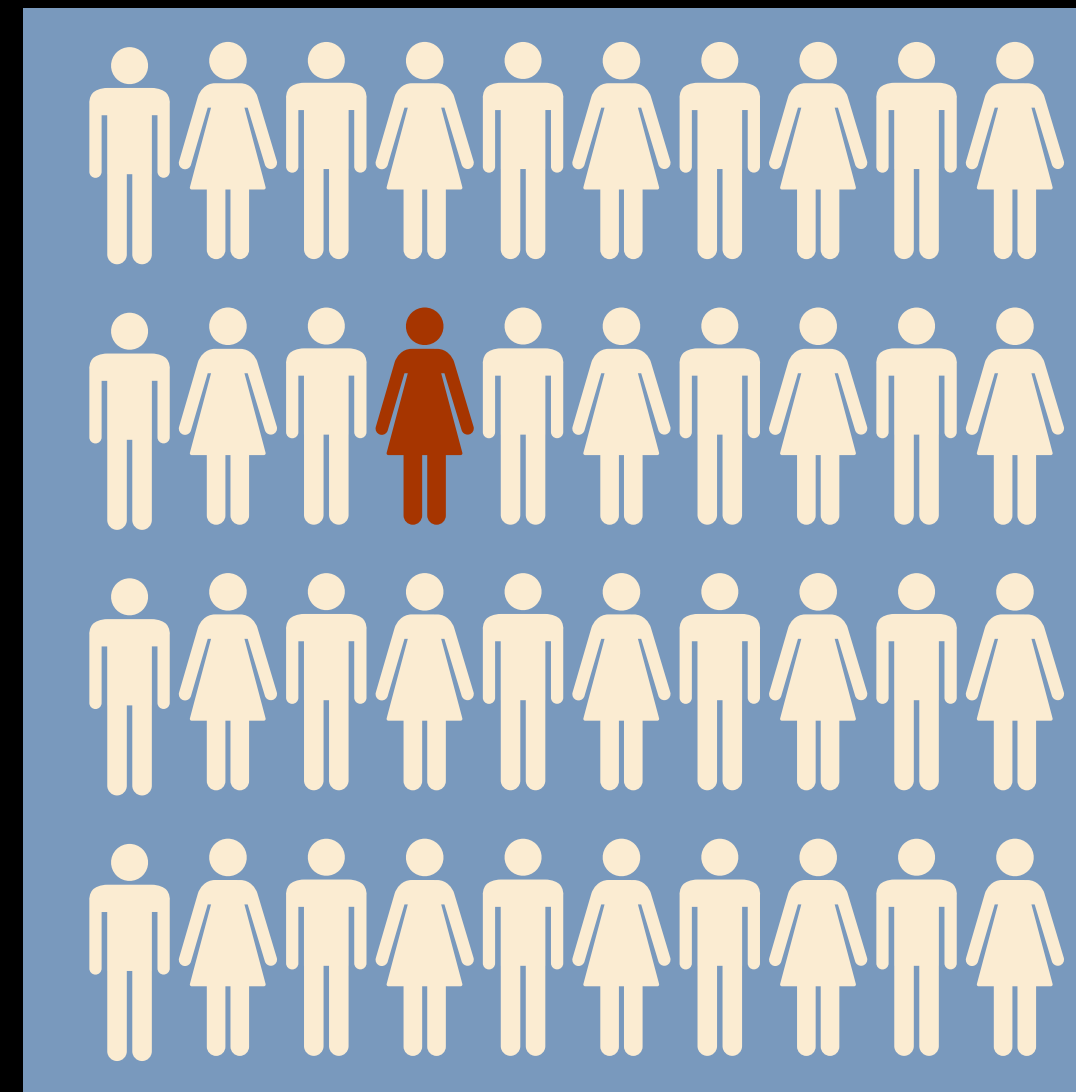
# Efficacy of a screening program

Fundamental challenge is **cost-effectiveness**

How much harm does the program do?

How much benefit does it achieve?



NLST Elligble　　Not

PMID: 23060474

1000 screens

↓

240 positives

↓

6 cancers

# Summary

All screening programs are classifiers

Effective screening programs need risk models to allocate care

Logistic Regression: Log-linear hypothesis class

Optimization: (Stochastic) Gradient Descent

Model selection and evaluation

# Questions?