

## Lecture 2

# Supervised Learning

## Part 1

TAs

Ajinkya Deogade

Amrita Singh

Diptodip Deb

Nils Eckstein

Virginia Rutter

Srini Turaga

turagas@janelia-hhmi.org

### 1. Key elements of SL/ML

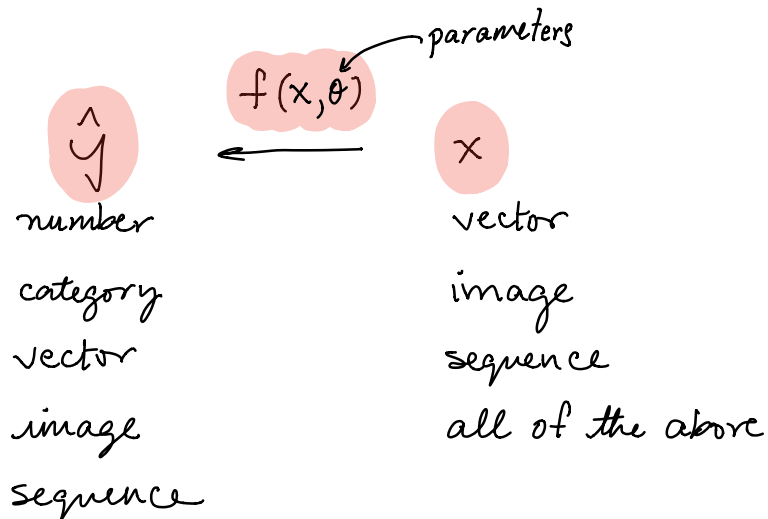
- data "x", "y"
- prediction function  $\hat{y} = f(x, \theta)$
- loss function  $\mathcal{L}(y, \hat{y})$
- optimization / training algorithm
- regularizer (optional)  $R(\theta)$

### 2. Linear regression

### 3. Classification as generalized linear modeling

### 4. Loss functions can be derived from probabilistic interpretations

# The what and why of supervised learning



## Black-box vs Interpretable

# The what and how of supervised learning

1.  $\begin{matrix} x_1, y_1 \\ x_2, y_2 \\ x_3, y_3 \\ \vdots \\ x_N, y_N \end{matrix}$   
Examples/data
2.  $\hat{y}_i = f(x_i, \theta)$   
prediction function
3.  $\mathcal{L}(y_i, \hat{y}_i)$   
scoring/objective/loss function
4.  $R(\theta)$   
regularization function  
(optional)
5.  $\min_{\theta} \sum_{i=1}^N \mathcal{L}(y_i, \hat{y}_i) + R(\theta)$   
 $= \min_{\theta} \sum_{i=1}^N \mathcal{L}(y_i, f(x_i, \theta)) + R(\theta)$   
optimization algorithm

# Generalized Linear Models

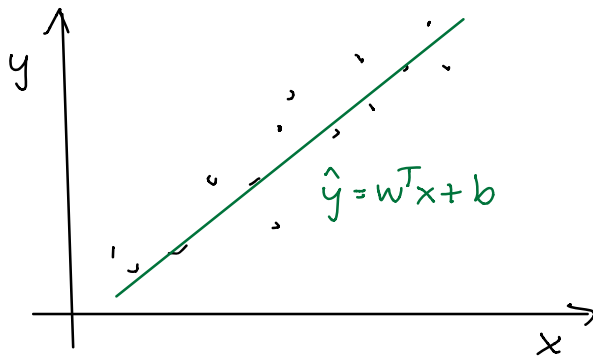
$$u = w^T x + b$$

$$\hat{y} = f(u)$$

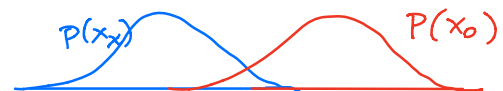
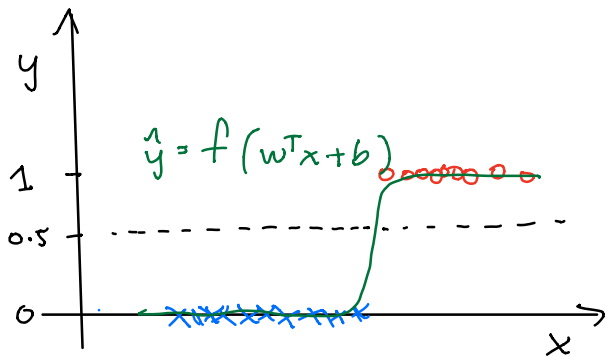
← linear  $\theta = \{w, b\}$

← nonlinear

Linear regression



Generalized Linear regression



Technical note:

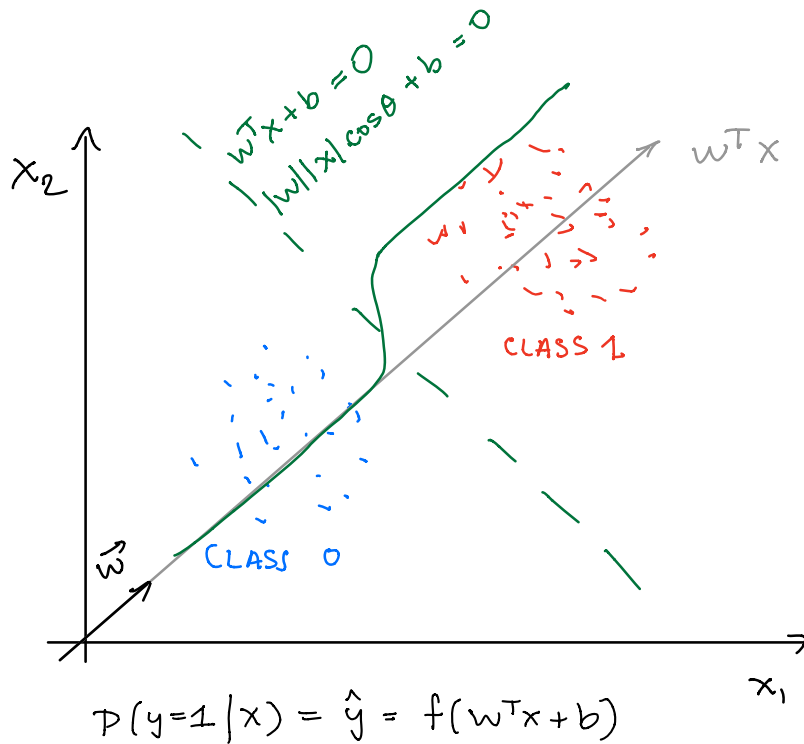
Equivalently  $\hat{y} = f(w^T x + b)$

↳  $f^{-1}(\hat{y}) = w^T x + b$

←  $f(u)$ : nonlinearity

←  $f^{-1}(y)$ : link function

# Logistic regression binary classification

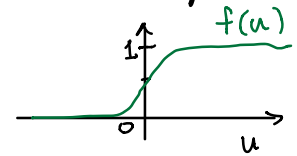


$$u = w^T x + b$$

$$\hat{y} = f(u)$$

$$f(u) = \frac{1}{1 + e^{-u}}$$

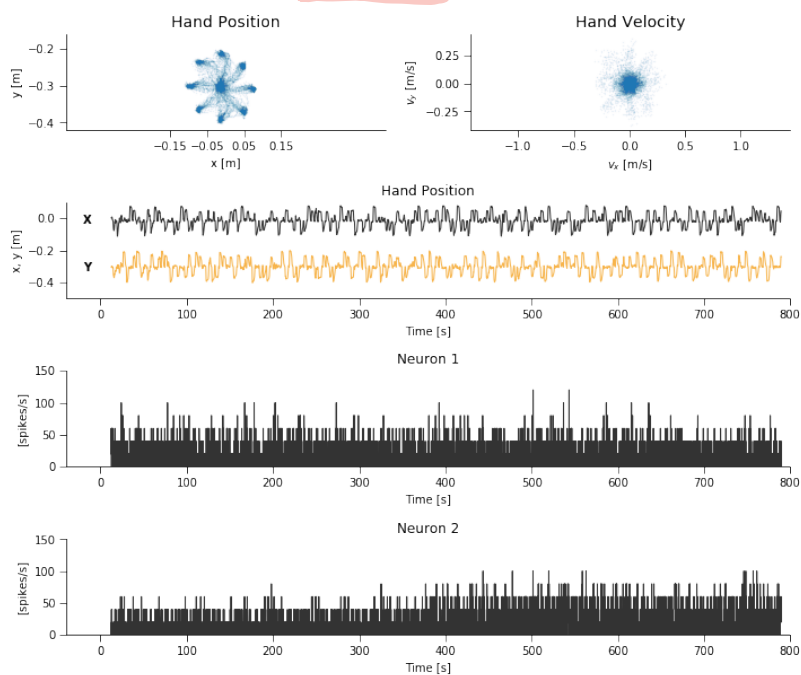
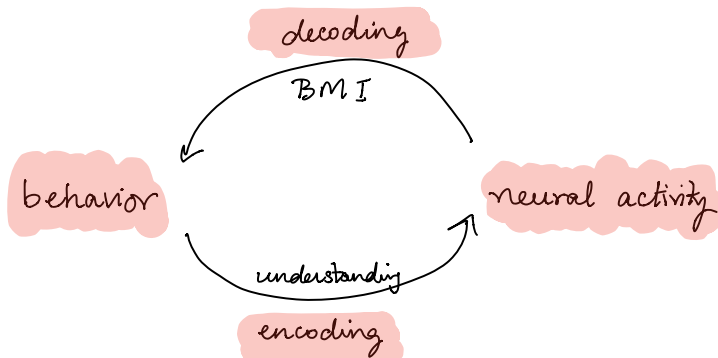
logistic sigmoid  
nonlinearity



$$\text{Loss function} : -\sum_i y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)$$

$$= -\log p(y_i | \hat{y}_i = f(w^T x + b))$$

# Neural data analysis



## Poisson regression

firing rate of a neuron at time  $t$

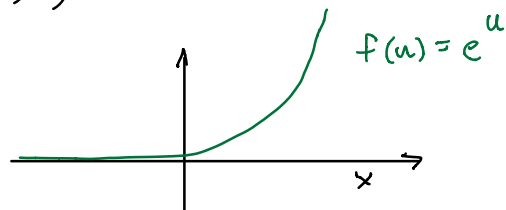
$$\hat{y}(t) = f(w^T x(t) + b) \quad \text{non-negative}$$

filter      stimulus

actual spikes are non-negative integers

$$y(t) \sim \text{Poisson}(\hat{y}(t))$$

So choose  $f(u) = e^u$



$$\text{Loss function: } \sum_i \hat{y}_i - y_i \log \hat{y}_i$$

$$= -\log P(y_i | \hat{y}_i = f(w^T x + b))$$

# Sparse regression

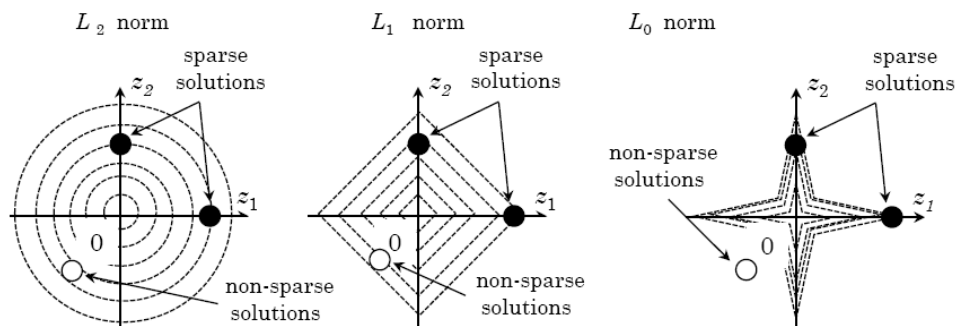
how to find informative features

$$\begin{aligned}\hat{y} &= w^T x + b \\ &= w_1 x_1 + w_2 x_2 + \dots + b\end{aligned}$$

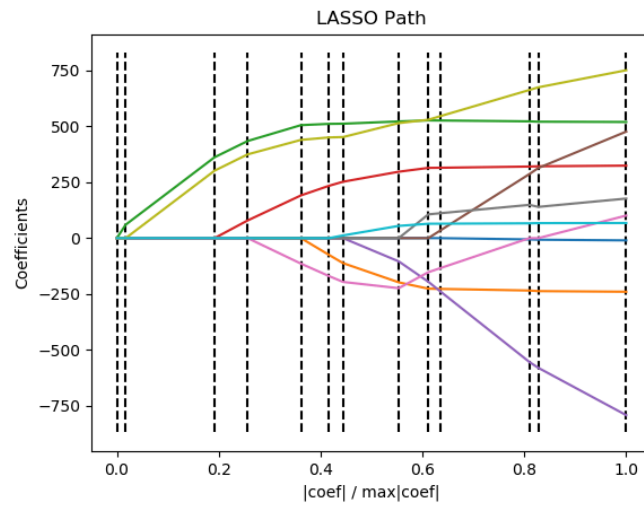
$$\begin{aligned}\text{Optimize : } \sum_i \mathcal{L}(y_i, \hat{y}_i) + \lambda \cdot R(w) & \text{ penalty for using} \\ & \text{ features} \\ &= \frac{1}{2} \sum_i (y_i - \hat{y}_i)^2 + \lambda \cdot R(w)\end{aligned}$$

Trade off between  $R(w)$  &  $\mathcal{L}(y, \hat{y})$

$$\text{Standard choice : } R(w) = \sum_j |w_j|^1 = \|w\|_1$$



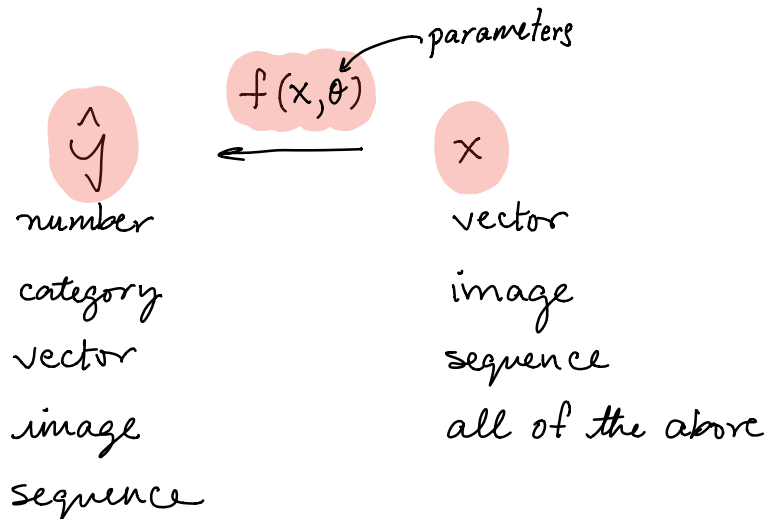
Lasso : sparse L1 regularized regression



← increasing  $\lambda$   
→ decreasing sparsity



## The what and why of supervised learning



## Black-box vs Interpretable

## The what and how of supervised learning

- $$\begin{matrix} x_1, y_1 \\ x_2, y_2 \\ x_3, y_3 \\ \vdots \\ x_N, y_N \end{matrix}$$
Examples/data
- $$\hat{y}_i = f(x_i, \theta)$$
prediction function
- $$\mathcal{L}(y_i, \hat{y}_i)$$
scoring/objective/loss function
- $$R(\theta)$$
regularization function  
(optional)
- $$\begin{aligned} & \min_{\theta} \sum_{i=1}^N \mathcal{L}(y_i, \hat{y}_i) + R(\theta) \\ &= \min_{\theta} \sum_{i=1}^N \mathcal{L}(y_i, f(x_i, \theta)) + R(\theta) \end{aligned}$$
optimization algorithm

