

Artificial Intelligence

CS 165A

Apr 7, 2022

Instructor: Prof. Yu-Xiang Wang

T
o
d
a
y

- Continuous optimization
- Wrapping up ML
- Starting PGM

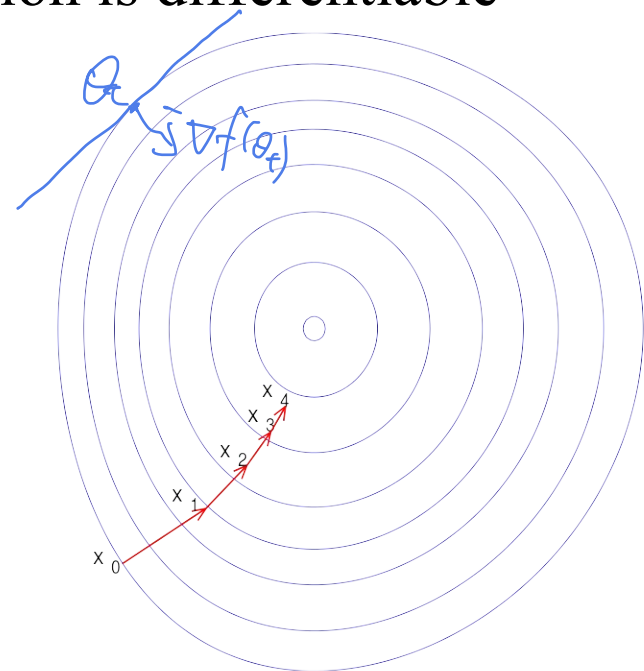
Recap: Last lecture

- Data splitting
 - Holdout and Cross validation for tuning hyperparameters
- The problem of distribution shift
- Learning a classifier:
 - From 0–1 loss to surrogate losses
 - Logistic loss as an example
- Continuous Optimization with Gradient Descent

Recap: How do we optimize a continuously differentiable function in general?

- The problem: $\min_{\theta} f(\theta)$
- Let's just optimize it anyway!
 - With gradient descent.
- Assumption: The objective function is differentiable almost everywhere.

$$\theta_{t+1} = \theta_t - \underbrace{\eta_t}_{\text{learning rate}} \underbrace{\nabla f(\theta_t)}_{\text{gradient}}$$



Recap: Gradient of logistic loss for learning a linear classifier

$$\nabla f(w) = \frac{1}{n} \sum_{i=1}^n \frac{\exp(-y_i \cdot x_i^T w)}{1 + \exp(-y_i \cdot x_i^T w)} (-y_i x_i)$$

Handwritten annotations: $w \in \mathbb{R}^d$, $x \in \mathbb{R}^d$. The entire expression is circled in blue. Arrows point from the handwritten text to the variables w and x_i in the formula.

- Question: What is the time complexity of computing this gradient?

Ans: O(nd)

Recap: Stochastic Gradient Descent (Robbins-Monro 1951)

- Gradient descent

$$\theta_{t+1} = \theta_t - \eta_t \nabla f(\theta_t)$$

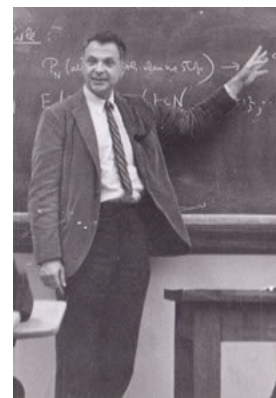
- Stochastic gradient descent

$$\theta_{t+1} = \theta_t - \eta_t \hat{\nabla} f(\theta_t)$$

- Using a stochastic approximation of the gradient:

$$\mathbb{E}[\hat{\nabla} f(\theta_t) | \theta_t] = \nabla f(\theta_t)$$

$$\mathbf{Var}[\hat{\nabla} f(\theta_t) | \theta_t] \leq \sigma^2$$



Herbert Robbins
1915 - 2001

The time complexity of each iteration in SGD can be drastically lower!

Plan for today

- Stochastic Gradient Descent
- Making sense of the SGD updates
- Multiclass classification
- Generative vs Discriminative Models
- Probability notations

One natural stochastic gradient to consider in machine learning

- Recall that

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(\theta, (x_i, y_i)) =: \mathcal{L}(\theta)$$

- Pick a **single** data point i uniformly at random

– Use $\nabla_{\theta} \ell(\theta, (x_i, y_i))$

$i \sim \text{Uniform} \{1, 2, \dots, n\}$

- Show that this is an unbiased estimator!

$$\mathbb{E}[\nabla_{\theta} \ell_i(\theta)] = \sum_{i=1}^n \frac{P(i)}{n} \cdot \nabla_{\theta} \ell_i(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell_i(\theta) = \nabla_{\theta} \mathcal{L}(\theta)$$

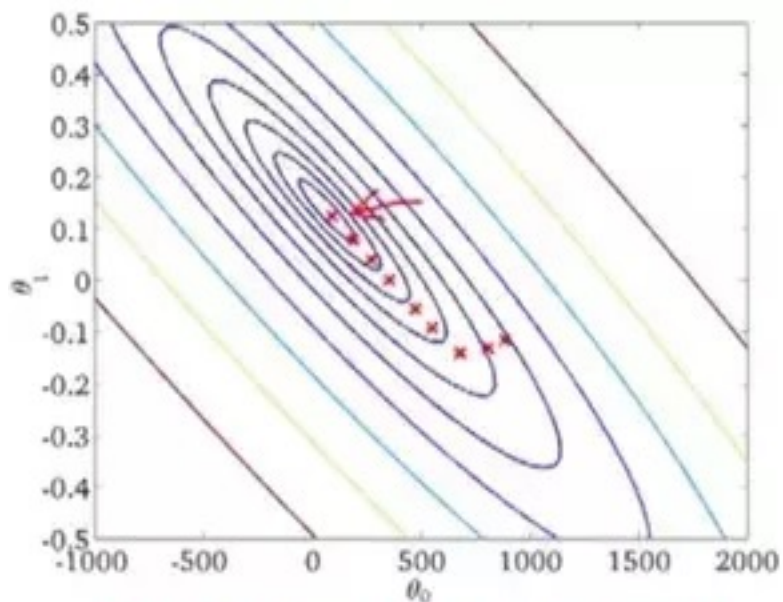
$$\nabla_{\theta} \ell_i(\theta, w) = \frac{\partial}{\partial w} \ell_i(\theta, w) \cdot X$$

Complexity: $O(d)$

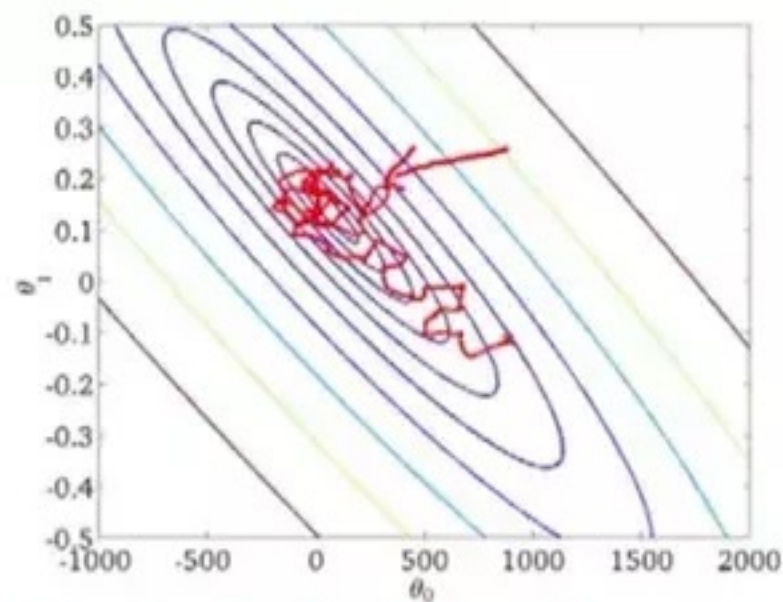
$$\downarrow$$

$$O(nnz(x))$$

Illustration of GD vs SGD



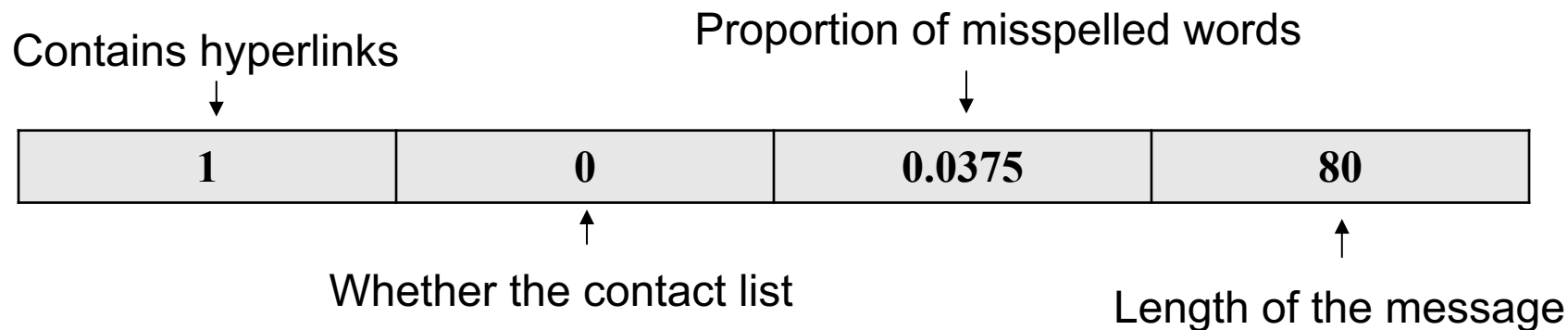
Batch Gradient Descent



Stochastic Gradient Descent

Observation: With the time gradient descent taking one step. SGD would have already moved many steps.

Intuition of the SGD algorithm on the “Spam Filter” example



- $\text{Score}(x) = w_0 + w_1 * 1(\text{hyperlinks}) + w_2 * 1(\text{contact list}) + w_3 * \text{misspelling} + w_4 * \text{length} = W^T x$ $w \in \mathbb{R}^5$

- **Meaning of these weights?**

- The more positive, the more we think the feature is associated with Spam email.
- The more negative, the less that we think the feature is associated with Spam email

Intuition of the SGD algorithm on the “Spam Filter” example

$$\nabla \ell(w, (x_i, y_i)) = \frac{\exp(-y_i \cdot x_i^T w)}{1 + \exp(-y_i \cdot x_i^T w)} \underbrace{\left(-y_i x_i \right)}_{\substack{\text{Scalar} > 0: \\ \approx 0 \text{ if the prediction} \\ \text{is correct} \\ \approx 1 \text{ otherwise}}} \begin{pmatrix} 1 \\ 0 \\ 0.375 \\ 80 \end{pmatrix}$$

$y=1$

Handwritten notes:

$$w_{t+1} = w_t + \eta \cdot \nabla \ell_i$$

$$\propto \begin{pmatrix} 1 \\ 0 \\ 0.375 \\ 80 \end{pmatrix}$$

Scalar > 0:
 ≈ 0 if the prediction
 is correct
 ≈ 1 otherwise

Vector of dimension d:
 provides the direction
 of the gradient

If we receive an example [1, 0, 0.0375, 80] like the one before.
 And a label y = 1 saying that this is a spam.

How does the SGD update change the weight vector?

Then by moving w towards the negative gradient direction, we are changing the weight vector by increasing the weights. i.e., increasing the amount they contribute to the score function (if currently the classifier is making a mistake on this example)

How to choose the step sizes / learning rates?

- In theory:

- Gradient decent: $1/L$ where L is the **Gradient Lipschitz constant** of the function we minimize.

- SGD: $\sum_t \eta_t = \infty, \sum_t \eta_t^2 < \infty$

- e.g. $\eta_t \in [1/t, 1/\sqrt{t})$

$$\sum_{t=1}^n \frac{1}{t} = 1 + \frac{1}{2} + \frac{1}{3} + \dots = O(\log n)$$
$$\sum_{t=1}^n \frac{1}{t^2} = 1 + \frac{1}{2^2} + \frac{1}{3^2} + \dots = \frac{1}{n^2} + \dots$$

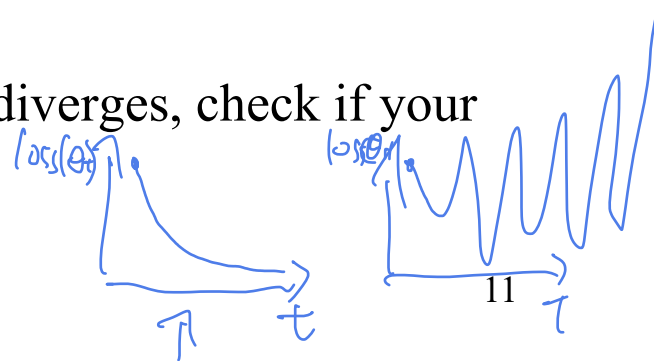
$$\approx \frac{\pi^2}{6} \approx 1.65 \dots \leq 2$$

- In practice:

- Use cross-validation on a subsample of the data.
- Fixed learning rate for SGD is usually fine.
- If it diverges, decrease the learning rate.
- If for extremely small learning rate, it still diverges, check if your gradient implementation is correct.

validation

"learning curve"



The power of SGD

- Extremely general:
 - Specify an end-to-end differentiable score function, e.g., a complex neural network.
 - Beyond the context of machine learning
- Extremely simple: a few lines of code.
- Extremely scalable
 - Just a few pass of the data, no need to store the data
- People are continuing to discover that many methods are special cases of SGD.

Multiclass classification problem

- You will be implementing a multi-class classification algorithm using linear logistic regression.

$$\mathcal{Y} = \{1, 2, \dots, k\}$$

- Recall that in binary classification you have one score function, the multi-class has one score function for each class.

$$h(x) = \underset{j \in \{1, \dots, k\}}{\operatorname{argmax}} \operatorname{Score}_{w_j}(x) = \underset{j \in \{1, \dots, k\}}{\operatorname{argmax}} (x^T \underline{w_j})$$

$$W = [w_1, w_2, \dots, w_k] \in \mathbb{R}^{d \times k}$$

- Now you have k linear score functions. Represent them as a matrix of weights W (**what's the dimension?**)
- You use the softmax-cross-entropy loss instead of the logistic loss.

Softmax Cross-Entropy Loss in the binary classification case is equivalent to the Logistic Loss we learned

- Softmax cross entropy loss

$$\text{loss}(W)(X, y) = - \sum_{i=1}^K \mathbb{1}(y=i) \log(\hat{p}_i^{(x)})$$

binary case: $-(y \cdot \log \hat{p} + (1-y) \cdot \log(1-\hat{p}))$

- In the binary case:

$$l(w) = \log(1 + \exp(-y \cdot \tilde{x}^T w))$$

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

when $y=1$, $l(w) = -\log(1) + \log(1 + \exp(-\tilde{x}^T w))$

$$= -(\log(1) - \log(1 + \exp(-\tilde{x}^T w)))$$

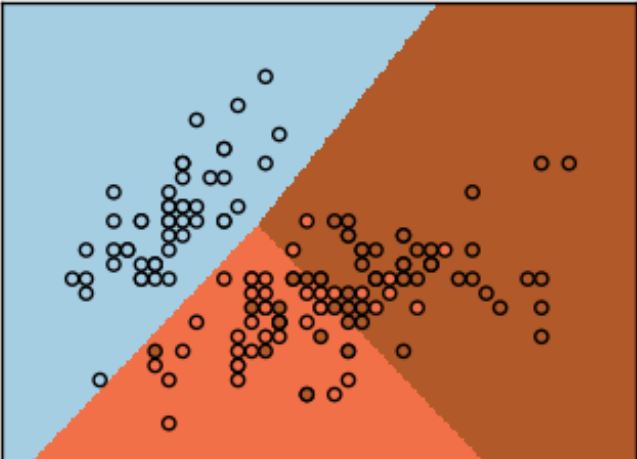
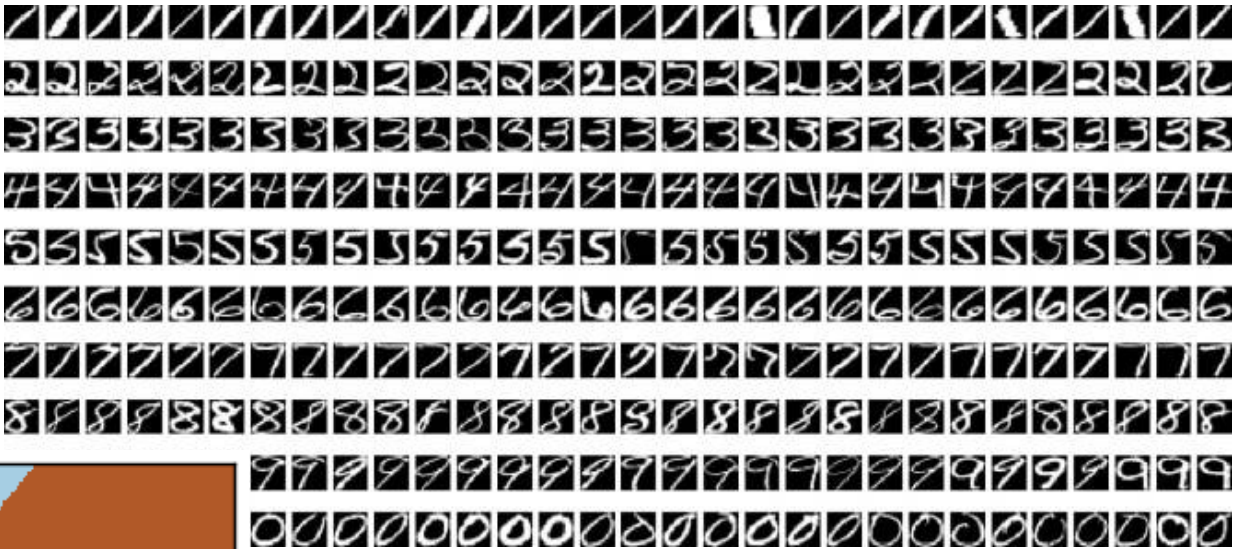
$$= -\log\left(\frac{1}{1 + \exp(-\tilde{x}^T w)}\right) = -\log\left(\frac{\exp(\tilde{x}^T w)}{1 + \exp(\tilde{x}^T w)}\right)$$

when $y=-1$. (exercise)

$$\begin{aligned} \hat{p}_i^{(x)} &= \text{softmax}_{\max}(\text{Score}(W; x)) \\ &= \frac{\exp(w_i^T x)}{\sum_{j=1}^K \exp(w_j^T x)} \end{aligned}$$

$$\boxed{\log x - \log y = \log \frac{x}{y} ?}$$

Decision boundaries of multi-class linear logistic regression



map image x to digit y

From linear logistic regression to neural networks

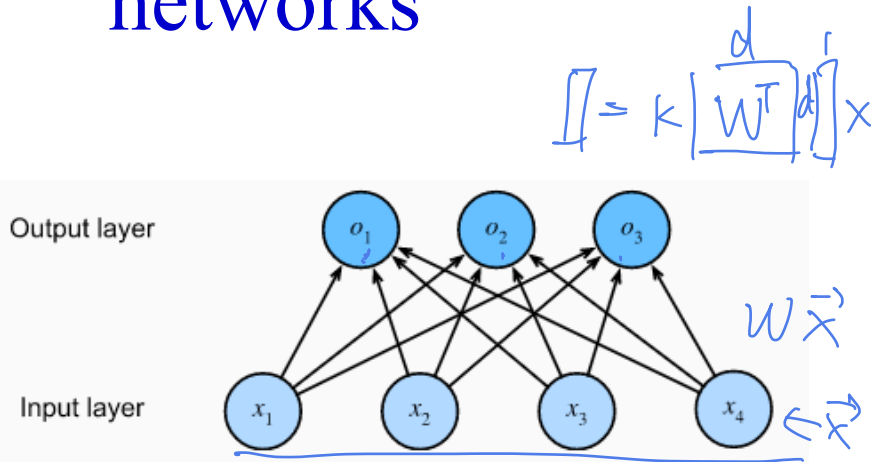


Fig. 3.4.1 Softmax regression is a single-layer neural network.

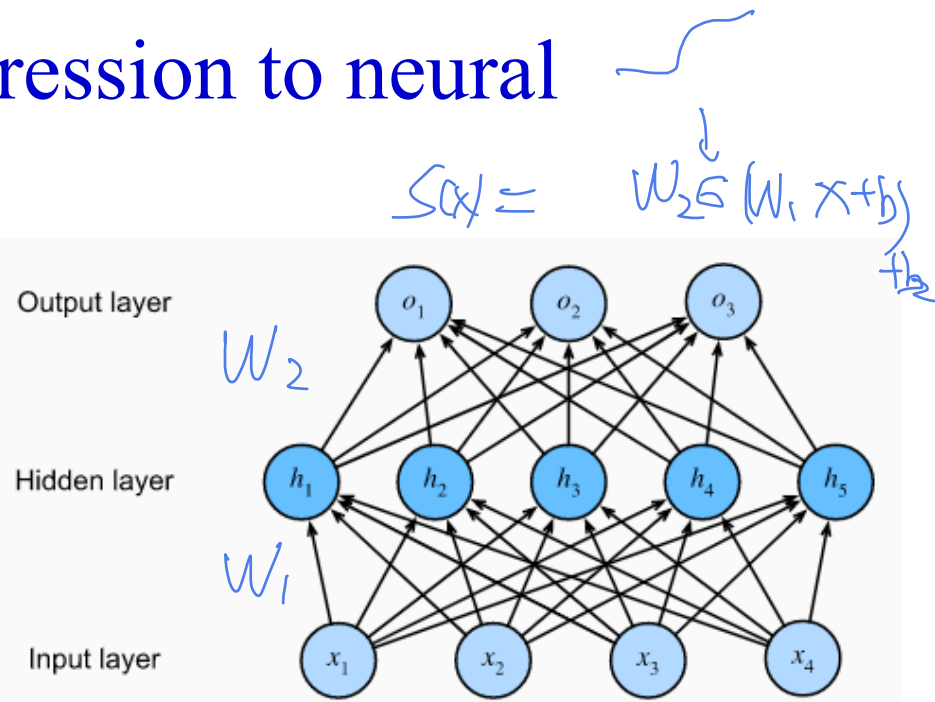


Fig. 4.1.1 An MLP with a hidden layer of 5 hidden units.

- Same input, same loss function
- The only difference is in how the score function is computed

From linear logistic regression to neural networks (side by side comparison)

Summary: steps to build a classifier agent

- 0. Collect a large labeled dataset.
- 1. Data splitting into training-validation-testing
- 2. Feature extraction
- 3. Specifying a hypothesis class
- 4. Learning with SGD using logistic loss
- 5. Validating on the validation set.
- 6. Testing on the test set.
- 7. Deploy the classifier agent.



Return to Step 2 or 3
If not happy with the
validation error.
Do model selection
(e.g. hyperparameter
tuning)

So far, we've learned everything about a classifier agent

- **Modeling**
 - Design meaningful feature extractors
 - Specify a family of classifiers on how the agent is going to classify examples using the features (indexed by free parameters)
- **Inference:** Trivially follows from the specification
- **Learning:** minimize errors (surrogate losses) over “free parameters”

This is called “discriminative” approach in modelling.

“Generative” modeling: Modeling the world with a (joint) probability distribution

- The “Discriminative approach” doesn’t care about the underlying processes that give rise to the observed data.
- **Modeling:** “Generative” modelling:
 - Label is a random variable Y , features are a vector random variables $X = [X_1, \dots, X_d]^T$
 - Model the world with a (family of) joint probability distributions.
- **Inference:** Then we can make principled inference by calculating $P(Y|X)$
 - You don’t just get a prediction, but also confidence.
- **Learning:** Find the distribution that fits the data well.

Side-by-side comparisons of two modelling principles: “Discriminative” vs “Generative”

	“Discriminative”	“Generative”
Modelling	<ul style="list-style-type: none">• Extract features• Hypothesis class• $h_{\theta}: X \rightarrow Y$	<ul style="list-style-type: none">• How data are generated• Family of probability distributions $P_{\theta}(X, Y)$
Inference	Just apply the classifier: $h(X)$	Calculate $P(Y X)$ Predict with: $\max_y P(Y=y X)$
Learning	Minimize error over “free parameters”: θ	Maximize likelihood over “free parameters”: θ

See Ng and Jordan (NIPS-2001):

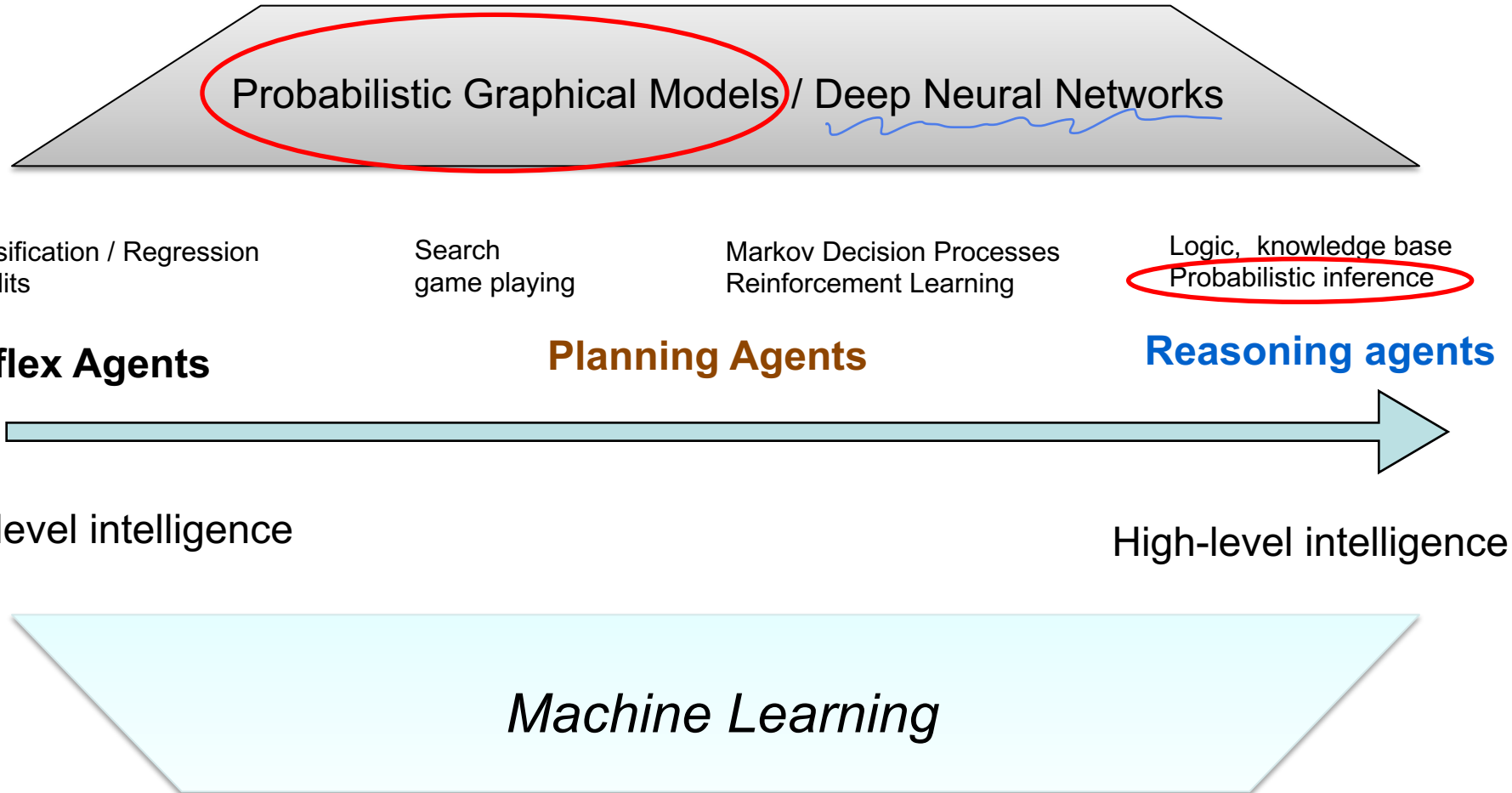
<https://ai.stanford.edu/~ang/papers/nips01-discriminativegenerative.pdf>

Example: Spam filter

- “Discriminative modeling”:
 - Think about features that has discriminative power
 - Think about hypothesis class (e.g., shape of the decision boundaries, functional form of the score function)
- “Generative modelling”:
 - Think about how “spammers” write (generate) their emails.
 - And also how “non-spammers” write their emails.
 - Model the probability of certain words appear / the probability of certain word combinations appear.

Once we learned probabilistic graphical model, we will be able to do this effectively.

Structure of the course



(Again this idea is adapted from Percy Liang's teachings)

Remaining time today and the next two lectures

- Probability notations
- Joint distributions, marginal, conditional
- Representing these quantities as arrays / matrices
- Conditional independences
- BayesNet, and examples
- d-separation, reasoning and inference, probabilistic modelling

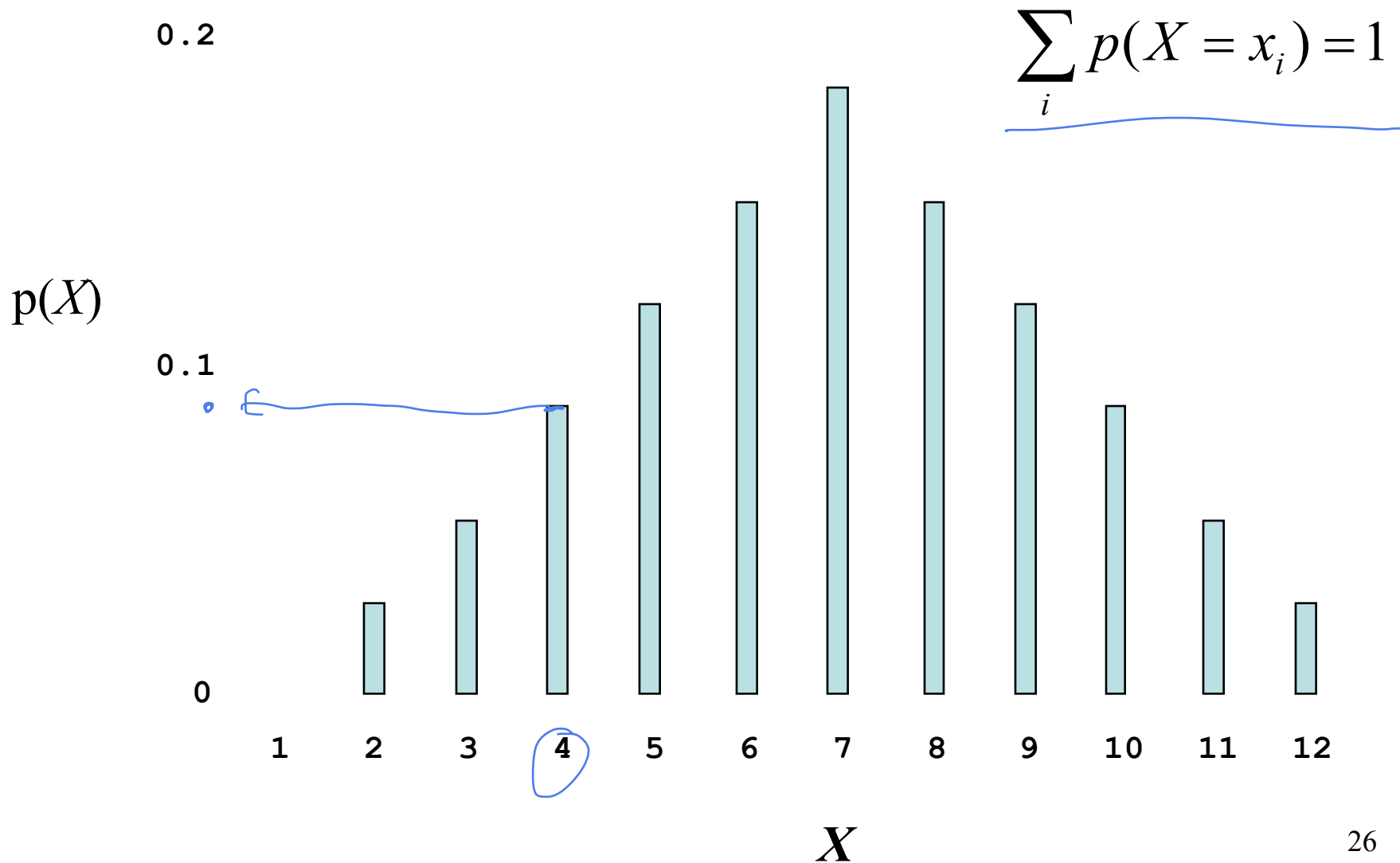
Probability notation and notes

- Probabilities of *propositions*
 - $P(A)$, $P(\text{the sun is shining})$
- Probabilities of *random variables (r.v.)*
 - $P(X = x_1)$, $P(Y = y_1)$, $P(x_1 < X < x_2)$
- $P(A)$ usually means $P(A = \text{True})$ (**A is a proposition, not a variable**)
 - This is a probability **value**
 - Technically, $P(A)$ is a probability *function*
- $P(X = x_1)$
 - This is a probability **value** ($P(X)$ is a probability *function*)
- $P(X)$
 - This is a **probability distribution** function, a.k.a. probability mass function (p.m.f.) for discrete r.v. or a probability density function (p.d.f.) for continuous r.v.
- Technically, if X is an r.v., we should not write $P(X) = 0.5$
 - But rather $P(X = x_1) = 0.5$

$$p(x) = P(X=x)$$

$P(X) = 0.5$
p.d.f.
p.m.f.

Discrete probability distribution



Continuous probability distribution

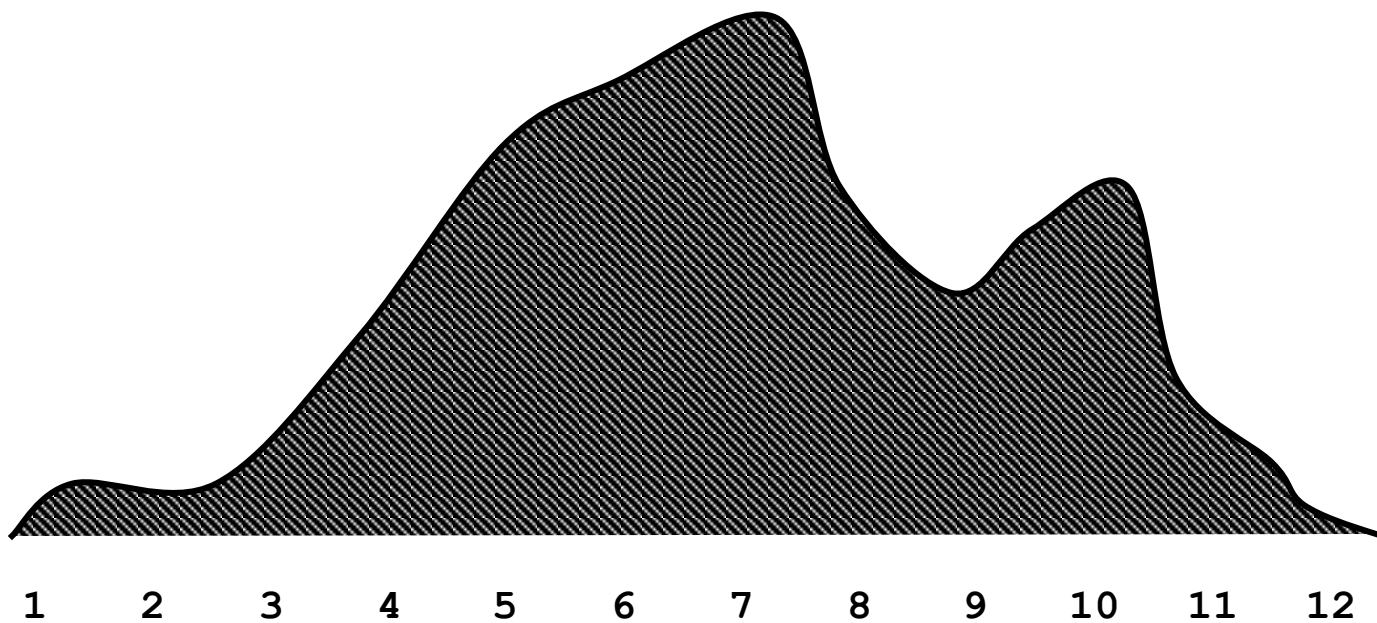
0.4

$$\int_{-\infty}^{\infty} p(X) = 1$$

$p(X)$

0.2

0



X

Joint Probabilities

- A **complete probability model** is a single joint probability distribution over all propositions/variables in the domain
 - $P(X_1, X_2, \dots, X_i, \dots)$
- A particular instance of the world has the probability
 - $P(X_1=x_1 \wedge X_2=x_2 \wedge \dots \wedge X_i=x_i \wedge \dots) = p$
- Rather than stating knowledge as
 - $\text{Raining} \Rightarrow \text{WetGrass}$
- We can state it as
 - $P(\text{Raining}, \text{WetGrass}) = 0.15$
 - $P(\text{Raining}, \neg\text{WetGrass}) = 0.01$
 - $P(\neg\text{Raining}, \text{WetGrass}) = 0.04$
 - $P(\neg\text{Raining}, \neg\text{WetGrass}) = 0.8$

	\neg WetGrass <small>not</small>	WetGrass
\neg Raining <small>not</small>	0.8	0.04
Raining	0.01	0.15

Marginal and Conditional Probability

- Marginal Probability
 - Marginal probability (distribution) of X : $P(X) = \sum_Y P(X, Y)$
 - **Bayesian interpretation:** Probabilities associated with one proposition or variable, **prior** to any evidence
 - E.g., $P(\text{WetGrass})$, $P(\neg\text{Raining})$
- Conditional Probability
 - $P(A | B)$ – “The probability of A given that we know B ”
 - **Bayesian interpretation:** After (**posterior** to) procuring evidence
 - E.g., $P(\text{WetGrass} | \text{Raining})$

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Assumes $P(Y)$ nonzero

or

$$P(X | Y) P(Y) = P(X, Y)$$

↑
Conditioned
margin

The chain rule: factorizing a joint distribution into marginal and conditionals

$$\underline{P(X, Y) = P(X | Y) P(Y)}$$

By the Chain Rule

$$\begin{aligned} P(X, Y, Z) &= \underline{P(X | Y, Z) P(Y, Z)} = P(X) P(Y | Z) \\ &= P(X | Y, Z) P(Y | Z) P(Z) && P(Z) P(Z | Y) \end{aligned}$$

or, equivalently

$$= P(X) P(Y | X) P(Z | X, Y)$$

- Notes:
- Precedence: ‘|’ is lowest
 - E.g., $P(X | Y, Z)$ means which?
 $P((X | Y), Z)$
 $P(X | (Y, Z)) \leftarrow$

Chain Rule implies Bayes' Rule

- Since $P(X, Y) = P(X | Y) P(Y)$

- and $P(X, Y) = P(Y | X) P(X)$

- Then $P(X | Y) P(Y) = P(Y | X) P(X)$

$$P(X | Y) = \frac{\overset{\text{conditional/likelihood}}{\downarrow} P(Y | X) \overset{\text{prior}}{\downarrow} P(X)}{P(Y) = \underbrace{\sum_X P(Y | X) P(X)}_{\text{normalization constant}}}$$

Bayes' Rule



Funny fact: Thomas Bayes is arguably a frequentist.

Stephen Fienberg. "When did Bayesian inference become 'Bayesian'?" *Bayesian analysis* 1.1 (2006): 1-40.
<https://projecteuclid.org/euclid.ba/1340371071>

Representing Probability Distributions using linear algebraic data structures (in python)

	<u>Continuous vars</u>	<u>Discrete vars</u>
<u>$P(X)$</u>	Function (of one variable)	<u>m vector</u>
$P(X=x)$	Scalar*	Scalar
$P(X,Y)$	Function of two variables	<u>m×n matrix</u>
$P(X Y)$	Function of two variables	m×n matrix
$P(X Y=y)$	Function of one variable	m vector
$P(X=x Y)$	Function of one variable	n vector
$P(X=x Y=y)$	Scalar*	Scalar

* - actually zero. Should be $P(x_1 < X < x_2)$

Example: Joint probability distribution

From $P(X, Y)$, we can always calculate:

- $P(X)$
 - $P(Y)$
 - $P(X|Y)$
 - $P(Y|X)$
- $P(X=x_1)$
 - $P(Y=y_2)$
 - $P(X|Y=y_1)$
 - $P(Y|X=x_1)$
 - $P(X=x_1|Y)$
 - etc.

		X		
		x_1	x_2	x_3
Y	y_1	0.2	0.1	0.1
	y_2	0.1	0.2	?

P(X,Y)

	x_1	x_2	x_3
y_1	0.2	0.1	0.1
y_2	0.1	0.2	0.3

P(X)

x_1	x_2	x_3
0.3	0.3	0.4

P(Y)

y_1	0.4
y_2	0.6

P(X|Y)

	x_1	x_2	x_3
y_1	0.5	0.25	0.25
y_2	0.167	0.333	0.5

$P(X=x_1, Y=y_2) = ?$

$P(X=x_1) = ?$

$P(Y=y_2) = ?$

$P(X|Y=y_1) = ?$

$P(X=x_1|Y) = ?$

P(Y|X)

	x_1	x_2	x_3
y_1	0.667	0.333	0.25
y_2	0.333	0.667	0.75

Quick checkpoint

- Probability notations
 - $P(A)$ is a number when A is an event / predicate.
 - $P(X)$ is a vector/function when X is a random variable.
- Joint probability distribution
 - Enumerating all combinations of events.
 - All values the random variables can take.
 - Assign a non-negative value to each.
- Marginals, conditionals
 - How they are related: Chain rule, Bayes rule