

3.2 connectionism: 1980s-1990s.

- ① idea: a large number of simple computational units can achieve intelligent behavior.

② distributed representation: each input should be represented by many features, and each feature should be involved in the representation of many possible inputs.

③ representation learning: learn representation

④ deep learning:

learn a hierarchy of representation.
i.e. learn right representation, and learn composition of representation.

3. history.

3.1 cybernetics: 1940s-1960s

③ Since 1995, kernel machines, graphical models led to decline in the popularity of deep learning.

3.3. deep learning: 2010 - today
success reason: have resources
deep learning can succeed.

- ① "Big Data" gives more data.
- ② hardware, software allow larger model size.

Chapter One

1. AI.

- ① want to create machines that think
- ② early days of AI tackled problems that can be described by mathematical rules.

③ true challenge to AI is solving tasks that we solve intuitively.
④ solution is allowing system to learn from experience and understand world in terms of a hierarchy of concepts.

2. ① knowledge base: hard-coded knowledge

- ② machine learning:
system has ability to acquire their own knowledge.

5.4. A is real symmetric

$\Rightarrow A$ has eigen decomposition
 $A = Q \Lambda Q^T$

Q is orthogonal matrix composed of eigenvectors of A .

Λ is diagonal, Λ_{ii} is eigenvalue:
 if $\lambda_i > 0, \forall i$, then A is called positive definite.

5.5. $A, \{x_i\}$ are eigenvalue of A .

$\Leftrightarrow \|x_i\|_2 = 1$
 $x_i^T x_j = 0$

4.2 vector x, y are orthogonal to each other

4.3 orthonormal of vectors
 \Leftrightarrow they are unit vector, and orthogonal to each other?

6. Singular Value Decomposition

$\forall A$ is matrix $\Rightarrow A$ has SVD

$\nexists A$ has eigen decomposition.
 $\Rightarrow A^T A = A A^T = I$

6.2 $A = U D V^T$ is called SVD of A .

U, V are orthogonal,
 $D \in \mathbb{R}^{m \times n}$ is diagonal.

this expression tells us some property about $1/2$.

5.2. v is eigenvector of A :
 $\det(A) := \lambda_1 \cdots \lambda_n$.

8. we use decomposition to represent A in other form,
 so we can see some different things. $\Leftrightarrow A = V \text{diag}(\lambda) V^{-1}$

Frobenius norm:

$$\|A\|_F = \sqrt{\sum_{i,j} A_{ij}^2}$$

4.1 unit vector x

$$\Leftrightarrow \|x\|_2 = 1$$

4.2 vector x, y are orthogonal

$$\Leftrightarrow x^T y = 0$$

4.3 orthonormal of vectors
 \Leftrightarrow they are unit vector, and orthogonal to each other?

5.1. $1/2 = 2 \times 2 \times 3$.

5.2. v is called a norm, if

$f(x) \geq 0, \forall x$

$f(x) = 0 \Rightarrow x = 0$

$f(x+y) \leq f(x) + f(y)$

$\forall x \in \mathbb{R}, f(\alpha x) = |\alpha| f(x)$

some property about $1/2$.

5.3. eigen decomposition of A

$\Leftrightarrow v \neq 0, Av = \lambda v$.

5.4. we use decomposition to

represent A in other form,

so we can see some different things. $\Leftrightarrow A = V \text{diag}(\lambda) V^{-1}$

Chapter 2
 linear algebra

2.3 $\sum_i c_i v_i$ is called linear combination of set $\{v_i\}$

2.4 set of linear combination of set V

is called span of V

2.5 a set $V = \{v^i\}, i=1, \dots, n\}$

is linear independent

\Leftrightarrow no vector in V is a linear combination of the other vectors.

1.2 $g(x, y) := f(y, x)$

g is called transpose of matrix f .

1.3 matrix product

$C = AB$, where $C_{ij} = \sum_k A_{ik} B_{kj}$

1.4 Hadamard product

$C = A \otimes B$, where $C_{ij} = A_{ij} B_{ij}$

2. $A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$

2.1 $Ax = b$ has solution

$\Leftrightarrow b$ in column space of A .

or say range of A .

3.1 L^p norm: $\|x\|_p = (\sum_i |x_i|^p)^{\frac{1}{p}}$

1.0 norm: $\|x\|_\infty := \max_i |x_i|$

2.2 $\forall b \in \mathbb{R}^m, Ax = b$ has solution

$\Leftrightarrow A$ has m linear independent columns.

1. definition

$f(x_1, \dots, x_n)$, where $x_i \in \mathbb{Z}^+$.

we write f in a table-like form

9. Common functions.

8. common pr. d.

9.1 logistic sigmoid

$$S(x) = \frac{1}{1+e^{-x}}, \forall x \in \mathbb{R}$$

9.2 softmax function

$$S(x) = \log(1 + e^x), \forall x \in \mathbb{R}$$

6. X and Y are independent, $X \perp Y \Leftrightarrow P(X=x, Y=y) = P(X=x)P(Y=y)$

8.1 Bernoulli distribution is a distribution over a single binary random variable. It's controlled by $\phi \in [0, 1]$.

8.2 multinomial or categorical ds. is a ds. over a single discrete variable with K different states, where K is finite.

11. Information theory

11.1 self-information of an event $X=x$ is $I(x) = -\log P(x)$

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

$$P(X|Y) = \lambda \log e^{XY}, \forall x, Y$$

11.2 Shannon entropy

$$H(P) = H(X) = \mathbb{E}_{X \sim P}[I(x)]$$

11.3 KL divergence

$$D_{KL}(P||Q) = \mathbb{E}_{X \sim P}[\log \frac{P(x)}{Q(x)}]$$

11.4 cross-entropy

$$H(P, Q) = H(P) + D_{KL}(P||Q)$$

11.5 Laplace dis.

$$p(x|\lambda) = \lambda \log e^{x\lambda}, \forall x \in \mathbb{R}$$

7.1 variance

$$\text{Var}(f(x)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$

7.2 standard deviation: $\sqrt{\text{Var}(f(x))}$

7.3 covariance

$$\text{Cov}(f(x), g(y)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])]$$

7.4 covariance matrix of a random vector $x \in \mathbb{R}^n$ is an $n \times n$ matrix.

$$\text{Laplace}(x; \mu, \gamma) = \frac{1}{2\pi} e^{-\frac{|x-\mu|}{\gamma}}$$

12. often the factorization of a probability ds. is represented with a graph, we call it a graphical model.

Chapter 3. Probability and Info theory

1. random variable X is a variable that can take on different values randomly.

2. for discrete variables, we describe probability distribution P using probability mass function (PMF),

$\Rightarrow P(X=x, Y=y) = P(X=x)P(Y=y)$
 $\forall x \in X, y \in Y$

for continuous variables, we use probability density function. (PDF).

3. marginal distribution is the probability distribution over a subset when we know the probability distr. over a set of variables

4. conditional probability is the pr. of some event given that some other event has happened.

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

5. chain rule or product rule of pr., which can get from definition of conditional pr., is

$$P(X^{(1)}, \dots, X^{(n)}) = P(X^{(1)}) \prod_{i=2}^n P(X^{(i)} | X^{(1)}, \dots, X^{(i-1)})$$

= $P(X^{(1)}) \prod_{i=2}^n P(X^{(i)} | X^{(1)}, \dots, X^{(i-1)})$

- 8.2. $f''(x_0) < 0$
 $\Rightarrow f'(x) \downarrow$ at x_0 .
 i.e. we downhill most (using 4.4).
 6.3. gradient descent:
 if $f(x)$ concave at x_0 , decrease f by moving in
 likewise, $f''(x_0) > f(x_0)$ convex at x_0 , in small steps with opposite sign of the derivative.
9. $f: \mathbb{R}^n \rightarrow \mathbb{R}$, then Hessian matrix is H , $H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$
- 9.2 directional second derivative
- 9.3 directional second derivative tells us how well we can expect a gradient descent step to perform.
- 9.4 $f(x_0) \approx f(x_0) + (x - x_0)^T H(x_0)(x - x_0)$ where g is gradient, H is Hessian at x_0 .
10. constrained optimization.
- if at least one feasible point exists, f is not permitted to have value ∞ , then $\min_{\text{S.t.}} f(x) \Leftrightarrow \max_{\lambda} L(x, \lambda, \alpha)$ where $S = \{g_i(x) = 0, h_j(x) \leq 0, \forall i, j\}$ where $L = f(x) + \frac{\alpha}{2} \sum_i g_i(x)^2 + \sum_j \lambda_j h_j(x)$
- [Chapter 4] 1. numerical computation
 1. on digital computer, we represent infinitely many real numbers with a finite number of bit patterns, incuring approximation error.
2. rounding error: underflow, overflow.
3. condition number of matrix A :

$$\frac{\max |A_{ij}|}{\min |A_{ij}|}, y = Ax, u = A^{-1}y$$
- when it is large, $f(x) = A^T x$ is sensitive to error in the input.
4. $f'(x_0)$ exists
 $\Leftrightarrow \forall \varepsilon > 0, \exists \delta > 0, \forall x \in D(x_0, \delta),$
 we have $\frac{f(x) - f(x_0)}{x - x_0} \in O(f'(x_0), \varepsilon)$
6. min $u^T \nabla f(x)$
 $\|u\|=1$
 $\Leftrightarrow \min_{\|u\|=1} \|u^T \nabla f(x)\| \text{ o.s.}$
7. Jacobian matrix.
 $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ then $J \in \mathbb{R}^{nm}$, $J_{ij} = \frac{\partial f_i}{\partial x_j}$
8. second derivative tells us how first derivative will change as we vary the input.
- 4.2 $f(x) - \alpha \text{sign}(f'(x)) < f(x)$ for small enough α ; $\alpha > 0$.
- 4.1. $f'(x_0) > 0 \Rightarrow \exists \delta > 0, \forall x \in D(x_0, \delta), f(x) \uparrow$ only at x_0 .
 $f'(x_0) < 0 \Rightarrow \exists \delta > 0, \forall x \in D(x_0, \delta), f(x) \downarrow$ at x_0 .

5. parametric models learn a

4.5 statistical learning theory

4. we can control whether a model is more likely to overfit or underfit by altering its capacity.

Averaged over all possible

data generating distributions, error is bounded from above by a quantity that grows as the model capacity grows but has the same error rate when classifying previously unobserved points.

averaging over all possible

data generating distributions, error is bounded from above by a quantity that grows as the model capacity grows but has the same error rate when classifying previously unobserved points.

averaging over all possible

data generating distributions, error is bounded from above by a quantity that grows as the model capacity grows but has the same error rate when classifying previously unobserved points.

averaging over all possible

data generating distributions, error is bounded from above by a quantity that grows as the model capacity grows but has the same error rate when classifying previously unobserved points.

averaging over all possible

data generating distributions, error is bounded from above by a quantity that grows as the model capacity grows but has the same error rate when classifying previously unobserved points.

averaging over all possible

data generating distributions, error is bounded from above by a quantity that grows as the model capacity grows but has the same error rate when classifying previously unobserved points.

averaging over all possible

data generating distributions, error is bounded from above by a quantity that grows as the model capacity grows but has the same error rate when classifying previously unobserved points.

Chapter 5, Machine learning basics.

1. learning algorithms algorithm is said to learn from experience with respect to task T and performance measure P , if P at T improves with E .

2. a machine learning algorithm specifies which family of functions the learning algorithm is an algorithm that is able to learn from data.

3. the central challenge in machine learning is that we must generalize well.

4.1 learning algorithm's effective capacity may be less than representational capacity of model, due to limitations of the capabilities of the optimization algorithm.

4.2 model specifies which family of functions the learning algorithm is an algorithm that is able to learn from data.

4.3 learning algorithm's effective capacity may be less than representational capacity of model, due to limitations of the capabilities of the optimization algorithm.

4.4 Occam's razor among competing hypotheses that explain known observations equally well, one should choose the simplest one.

5. we must design our machine algorithms to perform well on a specific task. (from 7).

- generating distribution, $P_{\text{Model}}(x; \theta)$ model distribution.
- 11.1. bias of an estimator $\hat{\theta}_m = E(\hat{\theta}_m) - \theta$
 - 12.1 maximizing the likelihood is equivalent to minimization where expectation is over the data.
 - 12.2. from maximum likelihood principle we can derive specific functions that are good estimators for different models.
 - 12.3 under appropriate conditions, maximum likelihood estimator is consistent.
 - 13.1 frequentist perspective : true parameter θ is fixed but unknown. point estimate $\hat{\theta}$ is a function of dataset. Here dataset is random, so $\hat{\theta}$ is r.v.
 - 13.2 Bayesian perspective : dataset is observed and is not random. true parameter θ is unknown or uncertain and thus is a r.v.
 14. The Bayesian uses probability to reflect degrees of certainty of state of knowledge.
9. Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.
- 10 hyperparameters are settings that we can use to control the behavior of a learning algorithm. The values of hyperparameters are not adapted by the learning algorithm itself.
- 10.2 validation set is used to "train" the hyperparameters.
- training error < validation error < generalization error.
- 10.3 cross-validation.
- repeating training and testing computation on different randomly chosen subsets or splits of the original dataset.
11. statistic $\hat{\theta}_m$ is a function of data.
- where X is training data, f_{Data} empirical distribution defined by X , P_{Model} true data distribution $\theta_m = g(X^{(1)}, \dots, X^{(n)})$

18. machine learning algorithm recipe.

dataset, cost function, optimization procedure, a model.

18.1 linear regression algorithm.

dataset $\{(X, y)\}$

$$\text{cost function: } J(w, b) = -\mathbb{E}_{X, Y \sim p_{\text{data}}} \log p_{\text{model}}(y|x)$$

$$\text{model: } p_{\text{model}}(y|x) = \mathcal{N}(y; w^T x + b, 1)$$

optimization algorithm: solve $\min J(w, b)$ using normal equations.

19.1 curse of dimensionality
machine learning problems become exceedingly difficult when the number of dimensions in the data is high.

19.2 The exponential advantages conferred by the use of deep, distributed representations counter the exponential challenges posed by the curse of dimensionality.

20. manifold learning algorithms assume that the data lies along a low-dimensional manifold.

15.4 SVM applies $\phi(x)$ to all inputs, then learns a linear model in the new transformed space. It does not provide probability, but only outputs a class identity.

$$f(x) = b + \sum_i \alpha_i k(x, x^{(i)})$$

15.5 kernel machines.

$$\begin{aligned} \theta_{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} p(\theta|X) \\ &= \underset{\theta}{\operatorname{argmax}} \log p(X|\theta) + \log p(\theta) \end{aligned}$$

The category of algorithms that employ the kernel trick is known as

kernel machines or kernel methods.

15.6 other supervised algorithms.
 k -nearest neighbors, decision tree.

16. unsupervised learning algorithms.

17. insight of stochastic gradient descent is that gradient is an expectation which may be approximated

15.2 linear regression corresponds to the family $p(y|x; \theta) = \mathcal{N}(y; \theta_0 + \theta_1 x, 1)$

15.3 logistic regression (used for classification)

y is binary variable.
estimated using a small set of samples.

14.2 full Bayesian inference
unlike maximum likelihood approach that make predictions using a point estimate of θ , the Bayesian approach is to make predictions using a full Bayesian posterior distribution over θ .

14.3 MAP estimate.

$$\begin{aligned} \theta_{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} p(\theta|X) \\ &= \underset{\theta}{\operatorname{argmax}} \log p(X|\theta) + \log p(\theta) \end{aligned}$$

where X is dataset, $p(\theta)$ is prior dist.

15. supervised learning algorithms.
15.1 most supervised learning algorithms are based on estimating $p(y|x)$.

15.2 linear regression

Chapter 6 MLP

- 4.2 If $f(x;\theta)$ is used to predict mean of Gaussian, the form of output layer, then maximum likelihood estimation choose activation functions, is equivalent to mean squared error.
- 4.3 one recurring theme throughout our design is that the gradient of the cost function must be large and predictable enough to serve as a good guide for the learning algorithm.
- 4.4 If we use a sufficiently powerful dnn, i.e., our model is able to represent any function from a wide class of functions, with this class being limited only by features such as continuity and boundedness rather than by having a specific parametric form, we can view cost function a being a functional rather than just a function.
- 3.1. model :
 the form of output layer,
 then maximum likelihood estimation choose activation functions,
 is equivalent to mean squared error.
2. function $f = f_2(f_1(x))$,
 $f_2(x) = w^T x$, $f_1(x) = \phi(x;\theta)$.
 f_1 is not linear.
 including depth of network,
 like kernel machines, we use a generic ϕ .
- 1.2 the goal of MLP is to approximate some function f^* .
 'cause ϕ is usually based only on the principle of local smoothness, it often generalizes poor.
- 1.3 when MLP is extended to include feedback connections, it is called rnn.
- 1.4 cnn is a specialized kind of MLP.
- 1.5 Names.
- 2.2 manually engineer ϕ .
 then we can use cross-entropy, then we can use training data and model distribution as having this form $f = f_2(f_1(x))$,
 f_1 is all hidden layers, and is learned using maximum likelihood.
- $J(\theta) = -\sum_{x,y} \log p_{\text{model}}(y|x)$
3. design decisions of deep learning algorithm :
- 4.1 specifying a model $p(y|x;\theta)$ choosing optimizer, automatically determines a cost function $J(\theta)$.
- 4.2 are hidden layers, $f^{(3)}$ is output layer.
- 4.3 x is called input layer.
- If $f^{(i)} = (f_1^{(i)}, f_2^{(i)})$, $f_1^{(i)}, f_2^{(i)}$ is vector-to-scalar function, called a neuron.

A maxout unit can learn a piecewise linear, convex function with up to k pieces.

It can be seen as learning the activation function itself rather than just the relationship between units.

- 5.6. we can represent dm as $f(x; \theta)$ which provides the parameter $w = f(x; \theta)$ for our model $p(y|x; \theta) = p(y|w)$. We predict y using our model $p(y|x; \theta)$, not dm $f(x; \theta)$.

6. Hidden units

6.5 ReLU is based on the principle that models are easier to optimize if their behavior is close to linear.

6.6 prior to introduction of ReLU, most neural networks used logistic sigmoid activation function $g(z) = \tanh(z)$,

$$g(z) = \sigma(z),$$

or hyperbolic tangent activation function $g(z) = \tanh(z)$

6.7 other hidden units.

RBF unit, Softplus, Hard tanh,

softmax, or cos.

- is then to provide some additional transformations from the features h to complete learning a full probability distribution

- 5.2 rectified linear units for Bernoulli output distribution.

$$p(y|x) = N(y; \hat{y}, I)$$

$$\hat{y} = W^T h + b$$

- 5.3 sigmoid units for Bernoulli output distribution.

- use the activation function $h = g(z) = \max(0, z)$,
 $z = W^T x + b$

y is binary variable.

- 5.4 softmax units for absolute value rectification, multinomial output distributions.

leaky ReLU,

parametric ReLU.

- $\hat{y}_i = P(y=i|x)$
- $$\hat{y}_i = \text{softmax}(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$
- $$z = W^T h + b$$

5.5

- where y is discrete variable, with n possible values for group i , $\{c_{i-1}, c_i, \dots, c_k\}$.

5.6

- we can model arbitrary probability distribution.

- 4.5 We can learn one conditional statistic of y given x , instead of learning a full probability distribution

- the task that dm must perform. $p(y|x; \theta)$.

- 4.6 different cost functions give different statistics.

- suppose f lies within the class we optimize, $f_{\text{opt}} = \arg \min_f E_{x, y \sim P_{\text{Data}}}[L(y, f(x))]$

- if we use mean squared error,
 $L(y, f(x)) = \|y - f(x)\|^2$

- then $f(x) = E_{y \sim P_{\text{Data}}(y|x)}[y]$

- $L(y, f(x)) = \|(y - f(x))\|_1$,
then $f(x) = \text{median of } P_{\text{Data}}(y|x)$

- 4.7 cross-entropy cost function is more popular than mse or mae.

5. given dm model $f(x) = f_2(f_1(x))$,

- here f_2 is output layer, $h = f_1(x)$ is

- all hidden layers, the role of f_2

8.4 chain rule of calculus. 8.2 bp algorithm allows the info from the cost to $x \in \mathbb{R}^m$, $y \in \mathbb{R}^n$, $z \in \mathbb{R}$, $y = g(x)$, $z = f(y)$.

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

8.5 back-propagation algorithm. It is an algorithm that computes the chain rule, is designed to reduce the number of common subexpressions without regard to memory. It performs on the order of one Jacobian product per node.

8.6 bp algorithm is only one approach to automatic differentiation that is concerned.

8.1 forward propagation. An operation is a function. The input x provides the initial info that then propagates up to hidden units at each layer and finally produces \hat{x} .

$\hat{x} = x$. \hat{x} is accompanied by a set of allowable operations.

8.2 bp algorithm. The info to compute derivatives with how to compute derivatives algorithmically. It is not the only way or the optimal way of computing gradient.

7.4 why we choose model that is deep?

It's because, first, we can get exponential advantage

(from 7.3), and second,

we have a statistical reason that choosing a deep model encodes a belief that the function we want to learn should involve composition of several simpler functions.

7.5 distributed representation network. The distributed representation network is given enough hidden units.

8.3 computational graph language.

8.1 forward propagation.

8.2 back-propagation.

$$z = xy.$$

8.3 The number of linear regions covered by a deep rectifier network with d input, depth L , and n units per hidden layer, is $O((\frac{d}{n})^{d(L-1)} n^d)$, i.e., exponential in the depth L .

7. most neural networks have their layers connected in a chain structure, like $f(g(f^m(x)))$, but they need not to.

7.2 universal approximation theorem.

A MLP with a linear output layer and at least one hidden layer with any "squashing" activation function can approximate any Borel measurable function from one finite-dimensional space to another with any desired non-zero amount of error, provided that the

composition of several simpler functions.

7.3 The number of linear regions covered by a deep rectifier network

with d input, depth L , and n units per hidden layer, is $O((\frac{d}{n})^{d(L-1)} n^d)$,

much than along v_j .

- we shrink those small eigenvalues; $\nabla_{\omega} \tilde{J}(\omega) = H(\omega - \omega^*)$
- H , to make ω small, at the same time, ensuring $\tilde{J}(\omega)$ won't increase much from $\tilde{J}(\omega^*)$.

④ For linear regression,

$$\tilde{J}(\omega) = (X\omega - y)^T(X\omega - y)$$

$$\tilde{J}(\omega) = J(\omega) + \frac{\alpha}{2} \|\omega\|^2$$

$$\text{let } \frac{\partial \tilde{J}}{\partial \omega} = 0, \Rightarrow \omega = (X^T X + \alpha I)^{-1} X^T y$$

$X^T X$ is proportional to the covariance matrix $\frac{1}{m} X^T X$, and diagonal entries of $\frac{1}{m} X^T X$ correspond to variance of each input feature,

L^2 regularization causes learning algorithm to "perceive" input X as having higher variance.

2.3 different choices for Ω

- can result in different solutions being preferred
- add weight decay to J , we get

$$\alpha \tilde{w} + H(\tilde{w} - w^*) = 0$$

where \tilde{w} is $\arg \min_w J(w) + \frac{\alpha}{2} \|w\|^2$, as weight decay, ridge regression, data generating process.

Let $H = Q \Lambda Q^T$ be eigendecomposition, Tikhonov regularization.

$$\tilde{J}(\omega) = J(\omega) + \frac{\alpha}{2} \|\omega\|_2^2$$

- $Q^T \tilde{w} = \text{diag}(\lambda_i) Q^T w^*$
- $\lambda_i > \alpha$, effect of regularization is small;

$\lambda_i \ll \alpha$, components of $\tilde{J}(\omega) = J(\omega) + \frac{\alpha}{2} \|\omega\|^2$

2. parameter norm penalty is a kind of regularization. We limit the capacity of models by adding a parameter norm penalty $C_2(\theta)$ to objective function J .

② let \tilde{J} approximate to J .

H is pd,

$\tilde{J}(\omega) = J(\omega^*) + (\omega - \omega^*)^T g(\omega^*)$

$\therefore \tilde{J}(\omega - \omega^*)^T H(\omega - \omega^*) + \frac{1}{2} (\omega - \omega^*)^T H(\omega - \omega^*)$

$\exists \gamma > 0, \forall \alpha \in (0, \gamma)$,

$$g(\omega^* + \alpha v_i) > \tilde{J}(\omega^* + \alpha v_i).$$

- If we move along v_i from ω^* , \tilde{J} will increase

Chapter 7. regularization for Ω .

- The goal of regularization is to take a model from overfitting regime to regime where the model family being trained matched the true data generating process.

1.2 In practical all scenarios, we use $\|\omega\|$ to stand for $\|\omega\|_2$. fitting model is a large model that has been regularized appropriately.

2. parameter norm penalty almost always do find that the best almost always do find that the best

1. The goal of regularization is to take a model from overfitting regime to regime where the model family being trained matched the true data generating process.

2.2 For neural networks, we penalize only weights w , and leave biases b unregularized.

$\therefore \omega^*$ is local minimum, H is pd (positive definite), $g = 0$.

6.2 If Ω is L^2 norm, then weights are constrained to lie in an L^2 ball, some for L^1 norm.

7. We can use regularization to solve underdetermined problems. For example,

pseudoinverse can be interpreted as stabilizing underdetermined problems using regularization.

8. Dataset augmentation schemes can reduce generalization, thus is a kind of regularization.

Example: apply transformation, inject noise.

$$= \arg \min_{\theta} J(\theta) + \alpha^* S(\theta),$$

let $\theta^* = \arg \max_{\theta} L(\theta, \alpha^*)$

we can see that this is exactly regularized training of minimizing J .

1. norm penalty \Rightarrow constrained optimization.

Here sparse refers to the fact that some parameters have an optimal value of zero.

1'. L^1 regularization has sparsity property, L^2 does not have.

6. Norm penalty as constrained optimization.

1. $\min_{\theta} J(\theta)$, subject to $S(\theta) < k$, \Leftrightarrow MAP Bayesian inference with isotropic Laplace prior.

$\Leftrightarrow \theta^* = \arg \max_{\theta} L(\theta, \alpha)$

where $L(\theta, \alpha) = J(\theta) + \alpha(S(\theta) - k)$
is generalized Lagrange function.

1. If we fix $\alpha = \alpha^*$

let $\theta^* = \arg \max_{\theta} L(\theta, \alpha^*)$

$= \arg \min_{\theta} J(\theta) + \alpha^* S(\theta)$, 5. difference between L^1 and L^2 .

we can see that this is

exactly regularized training

of minimizing J .

we suppose H is diagonal, $3.3 L^2$ regularization \Leftrightarrow MAP Bayesian inference with Gaussian prior.

let $\hat{w} = \arg \min_w \tilde{J}(w) + \alpha \|w\|_2$,
 $\hat{w}_i = \begin{cases} w_i^* - \frac{\alpha}{H_{ii}}, & \text{if } w_i^* \geq \frac{\alpha}{H_{ii}} \\ w_i^* + \frac{\alpha}{H_{ii}}, & \text{if } w_i^* \leq -\frac{\alpha}{H_{ii}} \\ 0, & \text{otherwise.} \end{cases}$

if $p(w) \sim N(w; 0, \frac{1}{\lambda} I^2)$

$$\begin{aligned} \theta_{\text{map}} &= \arg \max_{\theta} p(\theta | x) \\ &= \arg \max_{\theta} \log p(x | \theta) + \log p(\theta) \end{aligned}$$

if $p(w) \sim N(w; 0, \frac{1}{\lambda} I^2)$
then $\log p(w) = K + (-\frac{1}{2} \|w\|^2)$
proportional to $\lambda w^T w$ in L^2 .

where K doesn't depend on w ,
and does not affect learning process.

4. L^1 regularization.

we use quadratic Taylor polynomial to approximate $J(w)$

$$at w^* = \arg \min J(w).$$

$$\hat{J}(w) = J(w^*) + \frac{1}{2}(w - w^*)^T H(w - w^*)$$

15. Bagging, short for bootstrap aggregating, is a technique for reducing generalization error by combining several models. It's an example of a general strategy in ml, called model averaging. Techniques applying this strategy are known as ensemble methods.
13. parameter sharing is a way to express our prior knowledge about suitable values of model parameters.
- 13.2 Parameter sharing used by CNN is an example of how to effectively incorporate domain knowledge into network architecture.
- 15.1 The reason that model averaging works is that different models will usually not make all the same errors on the test set.
- 15.2. consider a set of k regression models. Suppose each model makes an error ε_i on every example, with variances $E[\varepsilon_i^2] = \sigma^2$, and covariances $E[\varepsilon_i \varepsilon_j] = \rho \sigma^2$,
- shares parameters with a discriminative model of $P(y|x)$, one can then trade-off the supervised criterion $-\log P(y|x)$ with the unsupervised $-\log P(x)$ or generative one $-\log P(x,y)$. The generative criterion then expresses a particular form of prior belief about solution to supervised learning problem.
14. Sparse Representation We can place a penalty on the activations of units, encouraging their activations to be sparse. This indirectly impose a complicated penalty on parameters.
- 14.2 norm penalty regularization of representation:
- $$J(\theta) = J(\theta) + \alpha \Omega(h)$$
- where $\alpha \in [0, \infty)$, h is representation.
9. Noise Robustness
- 9.2 Noise injection to input can be much more powerful than parameter norm penalty.
- 9.3 Noise applied to weights can be interpreted as a stochastic implementation of Bayesian inference over the weights, it can also be interpreted as equivalent to a form of regularization, encouraging stability of the function to be learned.
11. multi-task learning is a way to improve generalization by pooling the examples arising out of several tasks.
- 9.4 We can also inject noise at output target.
10. Early stopping is probably the most commonly used learning to provide regularization in deep learning.
- If we construct models in which a generative model can be viewed as a hyperparameter selection algorithm either $P(x)$ or $P(y|x)$

16.3 dropout is more effective than other standard computationally inexpensive regularizers, such as weight decay, sparse activity regularization.

Then, dropout training is:

It's very computationally cheap and does not significantly limit type of

model or training procedure.

17. we can reduce error rate on original f.i.d. test set via adversarial training.

18. Many ml algorithms aim to overcome the curse of dimensionality using manifold hypothesis.

Manifold hypothesis can be used as a kind of regularization.

- ① regularization.
- ② optimization.

19. central themes of machine learning by evaluating $p(y|x)$ in

15. boosting constructs an ensemble with higher capacity than individual models. In this technique, ensemble is not a regularization.

$$\min_{\theta} \min_{\mu} EJ(\theta, \mu)$$

16. Dropout

16.4. parameter sharing makes it possible to represent an exponential number of models with a tractable

amount of memory, and make most models do not have to be explicitly trained.

16.5. we can use

weight scaling inference

rule to approximate Perceptron rule to approximate Perceptron

one model.

16.2 Specifically, dropout

② optimization.

15. boosting constructs an ensemble with higher capacity than individual models. So, ensemble is not a regularization.

Then, dropout training is:

$$Var[\varepsilon] = E[\varepsilon^2] = \frac{1}{K} V + \frac{K-1}{K} C$$

Dropout provides a computationally inexpensive but powerful method of regularizing a broad

If errors are perfectly uncorrelated, it provides an inexpensive approximation to training and evaluating a boosted ensemble of exponentially many neural networks.

1'. If errors are perfectly correlated and $c=V$, $Var[\varepsilon]=V$.
If errors are perfectly uncorrelated, and $c=0$, $Var[\varepsilon]=\frac{V}{K}$.

1'. On average, ensemble will perform at least as well as any of its members, and if members make independent errors, it will perform significantly better than its members.

15.3 bagging uses K different datasets to get K different members. 15.4 model averaging is an extremely powerful and reliable method for reducing generalization error.

Chapter 8 Optimization for training deep models.

8. SGD : stochastic gradient descent
6.4. diff structure
6.5 vanishing and exploding gradient problem.

8.2 SGD is the most used algorithm for deep learning.
5.1 it follows the gradient of true generalization error so long as no examples are repeated. ① ml acts indirectly.
6.6 inexact gradient : approximate an intractable gradient.

8.3 while stopping criterion not met do:
sample minibatch $\{x^{(i)}, y^{(i)}\}$,
6.7 difficult global structure.

```
 $\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$ 
```

6.1 ill-conditioning : i.e., ill-conditioning of Hessian H .

the performance of an optimization algorithm is extremely difficult.

end while.

8.4 learning with SGD can 7. random functions behavior: and some being with sometimes be slow.

7.1 in low-dimensional space, local minima are common.

7.2 for SGD, it is $O(\frac{1}{k})$ in convex problem.
10. Gromer-Rao bound

states that generalization cannot decrease faster than $O(\frac{1}{k})$, k is iteration time.

Critical high cost are far points with more likely to be saddle points.

5. minibatch stochastic gradient descent
6.5 vanishing and exploding gradient problem.
8.2 SGD is the most used algorithm for deep learning.
6.6 inexact gradient : approximate an intractable gradient.

8.3 while stopping criterion not met do:
sample minibatch $\{x^{(i)}, y^{(i)}\}$,
6.7 difficult global structure.

6.1 ill-conditioning : i.e., ill-conditioning of Hessian H .

6.2 local minima

some arising from model identifiability problem,
and some being with high cost.

6.3 plateaus, saddle points and other flat regions:

In higher dim, local minima are rare and saddle points are more common.

6. challenges in optimization average.

1. learning differs from pure optimization
② cost function can be written as

$$\min_{\theta} \mathbb{E}_{X, Y \sim p_{\text{data}}} L(f(x; \theta), y) = \frac{1}{m} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$$

③ training algorithms do not halt at a local minimum.

2. erm : empirical risk minimization.

3. surrogate loss function

3.1 learning algorithm minimizes a surrogate loss function. It rarely uses erm.

3.2 early stopping criterion is based on true underlying loss function.

θ_{ML} = argmax_θ $\sum_i \log p_{\text{Model}}(x^{(i)}, y^{(i)}; \theta)$

⇒ min J(θ) = -E_{X, Y ∼ p_{\text{data}}} log p_{Model}(x, y; θ)

13.4 RMSProp

① Root Mean Square Propagation

② update rule:

$$g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i^m L(f(x_i^{(t)}; \theta), y_i^{(t)})$$

$$\begin{aligned} r &= \rho r + (1-\rho) g \\ \Delta \theta &= -\frac{\varepsilon}{\sqrt{r}} \theta \end{aligned}$$

$$\theta \leftarrow \theta + \Delta \theta$$

13.2 Adagrad algorithm.

axis-aligned

③ learning rate is a hyper-parameter that has a significant impact on model performance.

learning rate is a parameter that has a significant impact on model performance.

13.3 Adagrad algorithm

adaptive gradient

update rule:

$$g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i^m L(f(x_i^{(t)}; \theta), y_i^{(t)})$$

$$r \leftarrow r + g^2$$

$$\Delta \theta \leftarrow -\frac{\varepsilon}{\sqrt{r}} \theta$$

$$\theta \leftarrow \theta + \Delta \theta$$

③ if $r_i = t g_i^2$, then $r_i = t g_i^2$

learning rate of i -th para

$$r_i = \frac{\varepsilon}{\sqrt{t}} = \frac{\varepsilon}{\sqrt{1-p}} \cdot \frac{1}{g_i}$$

12.2. we almost always initialize all weights in the model to zero.

12.3. algorithms with adaptive learning rates.

12.4. Root Mean Square Propagation

12.5. Besides random initialization, we can also initialize model parameters need to break symmetry between different units.

12.6. we almost always initialize all weights in the model to zero.

12.7. It derives from a physical analogy.

12.8. momentum in physics is mass times velocity.

12.9. the initial distribution has a large effect on both optimization and generalization.

12.10. Large initial weights yield a large effect.

12.11. why it can accelerate?

12.12. If momentum algorithm always observes gradient g , say t times, then stepsize is $\frac{1-\alpha^t}{1-\alpha} \cdot \varepsilon \|g\|$.

12.13. terminal velocity is $\frac{\varepsilon g}{1-\alpha}$, $0 < \alpha < 1$.

12.14. optimization perspective suggests that weights should be large enough to propagate info successfully, but some regularization concerns encourage making them smaller.

12.15. Nesterov momentum.

12.16. initial parameters need to break symmetry between different units.

12.17. learning rate $\frac{\varepsilon}{\sqrt{t}}$ $< \frac{\varepsilon}{\sqrt{1-p}}$, parameters using ml.

15.4 Polyak averaging

It consists of averaging together several points in the trajectory through parameter space visited by an optimization algorithm.

15.4.1 BFGS is to approximate H^{-1} with M_t that is iteratively refined by low rank updates.

15.4.2 BFGS spends less time by refining each line search, but training simple models on simple tasks require store inverse Hessian approximation M from one step to the next than & ca. before confronting the challenge of training the desired model to perform the desired task.

15.5 It is more important to choose a model family that is easier to optimize than to use a powerful optimization algorithm.

15.7 continuation methods.

15.8 curriculum learning is a continuation method. It's based on the idea of planning a learning process to begin by learning simple concepts and progress to learning more complex concepts that depend on these simpler concepts simultaneously. $\theta_{t+1} = \theta_t + \varepsilon p_t$

14.4 BFOS

with M_t that is iteratively refined by low rank updates.

15.2 batch normalization is a method of adaptive reparametrization that acts to standardize only the mean and variance of each unit in order to stabilize learning,

but allows the relationships between units and the nonlinear conjugate gradients.

15.9 statistics of a single unit to change.

15.3 block coordinate descent refers to minimize with respect to a subset of variables simultaneously. $\theta_{t+1} = \theta_t + \varepsilon p_t$

② $J_2(\theta) := J(\theta) + (\theta - \theta_0)^T \nabla J(\theta_0)/3$. 5 Adam

$+ \frac{1}{2}(\theta - \theta_0)^T H(\theta - \theta_0)$ ① adaptive moments.

② It can be seen as a variant on the combination of RMSprop

and momentum.

15.5 pretraining is strategies that involve refining each line search, but

training simple models on simple tasks require store inverse Hessian approximation M from one step to the next than & ca. before confronting the challenge of training the desired model to perform the desired task.

15.6 It is more important to choose a model family that is easier to optimize than to use a powerful optimization algorithm.

15.7 continuation methods. $\theta = \theta_0 - [H + \alpha I]^{-1} \nabla_\theta J(\theta)$

15.8 curriculum learning is a continuation method. It's based on the idea of planning a learning process to begin by learning simple concepts and progress to learning more complex concepts that depend on these simpler concepts simultaneously. $\theta_{t+1} = \theta_t + \varepsilon p_t$

15.9 statistics of a single unit to change.

15.10 block coordinate descent refers to minimize with respect to a subset of variables simultaneously. $\theta_{t+1} = \theta_t + \varepsilon p_t$

15.11 Newton's method based on using a second-order Taylor series expansion to approximate $J(\theta)$ near some point

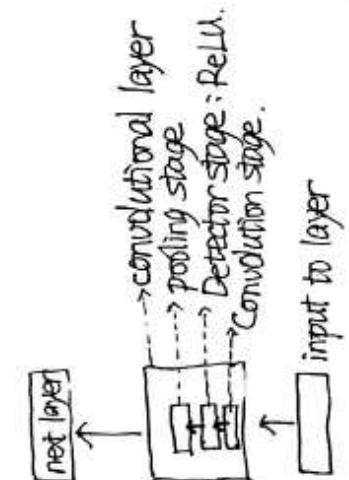
15.12 ignoring derivatives of higher order.

convolution \Rightarrow sharing parameter

3.7. any neural network that works with matrix multiplication and does not depend on specific properties of the matrix structure should work with convolution.

3.8. convolution is equivariant to translation.
 $f(x) \text{ is equivariant to } g(x)$
 $f(f(g(x)) = g(f(x))$.

3.9. a typical layer of CNN consists of three stages. See the figure.



- 3.3 the only reason to flip the kernel is to obtain the commutative property.
- 3.4 cross-correlation without any future changes to neural network.
- 3.5 we also call this operation convolution.
- 3.6 univariate discrete convolution corresponds to Toeplitz matrix.
- 3.7 convolution operation.
- 3.8 convolution is called convolution, block circulant matrix corresponds denoted as: $s(t) = (x * w)(t)$. where x is input, w is called kernel, s is called feature map.
- 3.9 convolution corresponds to a sparse matrix.
- 3.10 convolution is a linear operation. If w is a ptf, then s is a smoothed estimate of function x .
- 3.11 two-dimensional convolution
- 3.12 convolution \Rightarrow sparse interaction prior; convolution is a function of kernel, and the kernel is much smaller,
- 3.13 the matrix is sparse.

Chapter 9] Convolutional Networks.

1. convolutional networks are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers.

$$S(i,j) = (\mathbf{I} * K)(i,j)$$

1. a new best CNN architecture for a given benchmark is announced every few weeks to months.

$$S(t) = \sum_n x(n) w(t-n) da$$

this operation is called convolution, block circulant matrix corresponds denoted as: $s(t) = (x * w)(t)$.

where x is input, w is called kernel,

3.6 convolution corresponds to a sparse matrix.

3.7 convolution is a linear operation. If w is a ptf, then s is a smoothed estimate of function x .

3.8 convolution corresponds to a sparse matrix.

3.9 convolution is a function of kernel, and the kernel is much smaller,

3.10 convolution is a synonym for parameter sharing.

3.11 the matrix is sparse.

3.12 convolution is commutative.

3.13 convolution is commutative.

9.5 let K is a 4-D kernel tensor with $K_{i,l,m,n}$, where i indexing output channel, l input channel, m row offset within the input, n column offset within the input.

definition of convolution as used in mathematics.

9.3 convolution with a single kernel can only extract one kind of feature.

2 is produced by convolving K across V without flipping K .

8. The use of convolution is an infinitely strong prior that the function the layer should learn contains only local interactions and is equivariant to translation.

Likewise, the use of pooling is an infinitely strong prior that each unit should be invariant to small translations.

$$Z_{ijk} = \sum_{l,m,n} V_{i,j+m,k+n} K_{i,l,m,n}.$$

9.6. downsampled convolution function.

$$Z_{ijk} = c(K, V, s)_{i,j,k} \\ = \sum_{l,m,n} V_{i,j+sm,k+sn} K_{i,l,m,n}.$$

downsampled convolution is equivalent to convolution followed by downsampling, but has less computational cost.

9.7. zero padding the input allows us to control the kernel width and the size of output independently.

7.3. Pooling over spatial regions produces invariance to small translation.

If pooling over the outputs of separately parameterized convolutions, the features can learn which transformations to become invariant to.

7.4. pooling with downsampling can improve computational efficiency.

7.5. pooling is essential for handling inputs of varying size.

7.6. It is also possible to dynamically pool features together rather than using pooling function, for example, by running a clustering algorithm on the locations of interesting features.

8.3. some CNN architectures are designed to use pooling on some channels but not on other channels.

9. variants of basic convolution fourth axis indexing different examples in the batch.

9.2 convolution in CNN does not correspond precisely to the

13.1 one way to reduce the cost of convolutional network training is to use features that are not trained in a supervised fashion.

of the same kind of thing, for example, recordings over time, different widths of observations over space.

13.2 There are three ways for obtaining convolution kernels without supervised training. One is simply initialize them randomly. Another is to design them by hand. Finally, one can learn kernels with an unsupervised criterion.

13.3 Random filters often work surprisingly well in CNN.

13.3 Random filters often work surprisingly well in CNN. layers consisting of convolution following point-wise multiplication of the two signals, and converting back to the time domain using an inverse Fourier transform.

13.3 If kernel is separable,

First evaluate the performance of several convolutional network architectures by training only the last layer. Then take the best of these architectures. And train the entire architecture using a more expensive approach.

9.8. locally connected layer
(or unshared convolution)

$$Z_{l,j,k} = \sum_{l,m,n} V_{l,j,m,n} W_{i,j,k}$$

10.2 The model might emit a tensor S, where $S_{ij,k}$ is probability that pixel (j,k) belongs to class i .

12. efficient convolution algorithms. It is possible to speed up convolution of the input to the network by selecting an appropriate convolution algorithm.

12.2 convolution is equivalent to converting both the input and the kernel to the frequency domain instead of M.P. If not, using a Fourier transform, performing point-wise multiplication of the two signals, and converting back to the time domain using an inverse Fourier transform.

11. If computational cost and overfitting are significant issues, we can choose to use CNNs of kernel size k.

11.2 The use of convolution for processing variable sized inputs only makes sense for inputs that have variable size because they contain varying amounts of observation.

9.8. locally connected layer
(or unshared convolution)

$$Z_{l,j,k} = \sum_{l,m,n} V_{l,j,m,n} W_{i,j,k}$$

9.9 tiled convolution offers a compromise between a convolutional layer and a locally connected layer.

9.10 convolutional layers are hard-coded to be invariant to translation, while locally connected layers and tiled convolutional layers with max-pooling can make max-pooled units become invariant to some transformation.

10.1 CNNs can be used to output a tensor rather than just predicting a class label for a classification

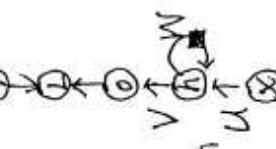
10.2 CNNs can be used to output a tensor rather than just predicting a class label for a classification

4.3 the network with recurrent connections only from the output at one time step to hidden units at the next time step is strictly less powerful because it lacks hidden-to-hidden recurrent connections.

4. models

universal Turing machine.

5. a computational graph is a way to formalize the structure of a set of computations.



6. The advantage of eliminating hidden-to-hidden recurrence is that, for any loss function based on comparing the prediction with the training target at time t to the training target at time t , all the time steps are decoupled. Training can thus be parallelized, with the gradient for each step t computed in isolation.

$$\text{gradient} = \frac{\partial L}{\partial w_1} + \frac{\partial L}{\partial w_2} + \dots + \frac{\partial L}{\partial w_n}$$

this RNN is universal in the sense that any function computable by a Turing machine can be computed by such a RNN of a finite size.

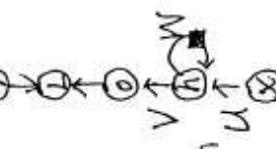
3.4 RNNs with recurrent connections between hidden units, that read an entire sequence and then produce a single output

$$4.2 \quad a(t) = b + W a(t-1) + U x(t)$$

for example, it cannot simulate a

① recurrent graph:

5. a computational graph is a way to formalize the structure of a set of computations.



6. The advantage of eliminating hidden-to-hidden recurrence is that, for any loss function based on comparing the prediction with the training target at time t to the training target at

①

$$\text{gradient} = \frac{\partial L}{\partial w_1} + \frac{\partial L}{\partial w_2} + \dots + \frac{\partial L}{\partial w_n}$$

with the gradient for each step t

computed in isolation.

this RNN is universal in the sense that any function computable by a

Turing machine can be computed by such a RNN of a finite size.

first, regardless of the sequence length, the learned model always has the same input size.

second, it is possible to use the same transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

some transition function f with the same parameters at every time step.

Chapter 10: Sequence Modeling

1. involves almost any function can be considered a MLP, essentially any function involving recurrence can be considered a RNN.

can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

be considered a RNN, essentially any function involving recurrence can be considered a RNN.

- 10¹) notes in the graphical model decouples the past and the future, and the model can be efficiently parametrized by using the same conditional probabilities distribution at each time step, and that, when the variables are all observed, the probability of the joint assignment of all variables can be evaluated efficiently.
- 10⁴ The price RNNs pay for their reduced number of parameters is that optimizing the parameters may be difficult.
- 10⁵. parameter sharing in RNNs \Leftrightarrow the conditional probability distribution over the variables at time $t+1$ given the variables at time t is stationary.

We parametrize the joint distribution of these observations using the chain rule for conditional probabilities:

$$P(Y) = P(y^{(0)}, \dots, y(t)) = \prod_{t=1}^T P(y(t)|y^{(t-1)}, \dots, y^{(0)})$$

RNNs are useful when we believe that the distribution over $y(t)$ may depend on a value of $y^{(i)}$ from the distant past in a way that is not captured by the effect of $y^{(i)}$ on $y^{(t-1)}$.

One way to interpret RNN as graphical model is to view the RNN as defining a graphical model whose structure is the complete graph, able to represent direct dependencies between any pair of y values.

10² consider the case where RNN models only a sequence of scalar random variables $Y = \{y^{(1)}, \dots, y^{(\tau)}\}$, with no additional inputs, x . The input at time step t is simply the output at time step $t-1$. Over the observations.

9. As with a MLP, we usually wish to interpret the output of the RNN as a probability distribution, and we usually use the cross-entropy loss criterion, in which during training the model receives the ground truth output $y^{(t)}$ as input at time $t+1$.
- 7.1 teacher forcing is a procedure associated with that distribution to define the loss. Mean squared error is the cross-entropy loss associated with an output distribution that is unit Gaussian, for example, just as with a MLP. 8. computing the gradient through a RNN simply applies the generalized back-propagation algorithm to the unrolled computational graph of RNN.
- 8.2 BPPT.
- 7.2 teacher forcing is a procedure that emerges from maximum likelihood criterion, in which during training the model receives the ground truth output $y^{(t)}$ as input at time $t+1$.
7. models that have recurrent connections from their outputs leading back into the model may be trained with teacher forcing.

The RNN then defines a directed graphical model over the y variables.

many stages tend to either vanish and associated transformations (most of the time) or explode (rarely). ① from input to hidden state. This problem is particular to RNNs. Various approaches have been proposed to reduce the difficulty of learning long-term dependencies.

18. Both ESNs and liquid state machines are termed reservoir computing to denote the fact that the hidden units form of reservoir of temporal features which may capture different aspects of the history of inputs.

18.2 ESNs: Echo State Networks.

18.3 The recurrent weights and input weights are some of the most difficult parameters to learn in a RNN. Reservoir computing avoids this difficulty by setting the recurrent

represent a distribution $P(y|x)$ by using the same $P(y|w)$ as

② from previous hidden state
before, but making w a function of x .
One option is adding an extra symbol to the vocabulary.

Another approach is to introduce an extra Bernoulli output to the model that represents the decision to either continue generation or halt generation at each time step. Another way to determine sequence length τ is to add extra output to another RNN that moves backward through time beginning from the end of the sequence. length τ is predicted by the model that predicts the integer τ itself.

14. the encoder-decoder or sequence-to-sequence architecture is a RNN architecture for RNNs, with a different kind of computational graph, which maps a variable-length sequence to another variable-length sequence.

15. the computation in most RNNs can be decomposed into three blocks of parameters

11. various ways to determine the length of sequence

③ from previous hidden state to the next hidden state.

13. bidirectional RNNs combine an RNN that moves forward through time beginning from the start of the sequence with another RNN that moves backward through time beginning from the end of the sequence. length τ is predicted by the model that predicts the integer τ itself.

12. modeling sequences conditioned on context with RNNs.
D.2 Any model representing a variable $P(y; \theta)$ can be reinterpreted as a model representing a conditional distribution $P(y|w)$ with $w = \theta$. We can extend such a model to

controls the forgetting factor and the decision to update the state unit.

- 21. optimization for long-term dependencies vanish nor explode.

21.2 gradient clipping helps to deal with exploding gradients, but it does not help with vanishing gradients. To address chosen constants or were parameters. **Leaky units** did this with connection weights that were either manually created paths in the computational graph of an unfolded recurrent architecture along with the product of gradients associated with arcs is near 1. One approach to achieve "cells" that have an internal recurrence mechanism. Another idea is to regularize the parameters so as to encourage "info flow".

22. neural networks excels at storing implicit knowledge. They struggle to memorize facts. Memory networks introduced to resolve this difficulty. They include a set of memory cells that can be thought of as an extension of the memory cells in LSTMs and GRUs.

It has three kinds of gating units, forget gate unit f , external input gate units g , and output gate q .

20.4 GRU : gated recurrent unit. The main difference with LSTM is the GRU is a single gating unit simultaneously gated RNNs, which include

idea of creating paths through time way of ensuring that the hidden units do a good job of computing the history of past inputs

unit can access values from the past.

19.3 hidden units with linear self-connections are called

19.4 another approach to handle long-term dependencies is the idea of organizing the state of the RNN at multiple timescale, with the flaying more easily through long distances at the slower time scale. It involves actively removes scales and transfer information from the distant past to the present more efficiently.

19.2 adding skip connections through time

Skip connections through d time steps are a way of ensuring that a unit can always learn to be influenced by a value from d time steps earlier.

The use of linear self-connection with a weight near one is a different gated RNNs are based on the

weights such that the recurrent hidden units do a good job of computing the history of past inputs

and learning only the output weights and learning only the output weights

19. one way to deal with long-term dependencies is to design a model that operates at multiple time scales, so that some parts of the model operate at fine-grained time scales and can handle small details, while other parts operate at coarse time

19.2 adding skip connections through time

Skip connections through d time steps are a way of ensuring that a unit can always learn to be influenced by a value from d time steps earlier.

The use of linear self-connection with a weight near one is a different gated RNNs are based on the

that is already known to perform best on the previously studied task, even copy a trained model from that task.

- (vi) add unsupervised learning later. It's possible to trade coverage 24 metrics for accuracy.
- 4. determine whether to gather more data.
 - (i) determine whether performance on the training set is acceptable.

If performance on the training set is poor, then try increasing the size of the model and try improving the learning algorithm, for example, by tuning the learning rate hyperparameter.

- examples for which the machine learning system is able to produce a response.
- 2.3 determine the target value of the performance metric and choose which metric to use.
 - (i) determine your goals: what error metric to use, and your target value for this error metric.
 - (i) accuracy \Leftrightarrow error rate
 - (ii) click-through rates, user satisfaction.
 - (iii) precision P: It is the fraction of detections reported by the model that were correct.
 - (iv) recall r: It is the fraction of true events that were detected.
- 3. default baseline models. After choosing performance metrics and goals, the next step is to establish a reasonable end-to-end system asap.
- 3.2 (i) choose the general category of model based on the structure of your data.
 - (ii) choose optimization algorithm such as SGD with momentum with a decaying learning rate, or Adam.
 - (iii) batch normalization
- (ii) if the performance on the training set is acceptable, then measure the performance on a test set.

Chapter 11] Practical Methodology

1. practical design process
 - (i) determine your goals: what error metric to use, and your target value for this error metric.
 - (ii) establish a working end-to-end pipeline as soon as possible.
 - (iii) instrument the system well to determine bottlenecks in performance.
 - (iv) repeatedly make incremental changes such as gathering new data, adjusting hyperparameters, or changing algorithms, based on specific findings from your instrumentation.
- 2.1 Bayes error defines the minimum error rate that you can hope to achieve, even if you can have infinite training data and can recover the true probability distribution.
- 2.2 F-score: $F = \frac{2PR}{P+R}$
- (i) early stopping, dropout
- (ii) coverage: It is the fraction of

1.2 Usually the best performance on the test set comes from a large model that is regularized well.

7.3 Your end goal is good performance on the test set. The brute force way to practically guarantee success is to continually increase model capacity and training set size until the task is solved. In principle, this approach could fail due to optimization difficulties.

8 debugging strategies

8.2 Some debugging tests:

- (i) visualize the model in action.
- (ii) visualize the worse mistakes.
- (iii) reasoning about software using training and test error.

low training error and high test error
⇒ model is overfitting or test error is measured incorrectly.
high training error and high test error
⇒ underfitting or software defect.
(iv) fit a tiny dataset.

high error on training set with few examples
⇒ software defect.

7.1 compare backpropagated derivatives to numerical derivatives.

(v) monitor histograms of activations and gradients. These provide guarantees about the quality of the optimization process.

wrap a learning algorithm and choose its hyperparameters, thus hiding the hyperparameters of the learning algorithm from the user.

6. selecting hyperparameters

If the performance on the test set is also acceptable, then there is nothing left to be done.

If the test set performance is much worse than the training set performance, then gathering more data is one of the most effective solutions.

If it is costly to gather more data, then reduce the size of the model or improve regularization.

The generalization error typically plotted as a function of one of the hyperparameters. If you have time to tune only one hyperparameter, tune the learning rate. The effective capacity of the model is then gathering more data is advisable.

Most model-based algorithms for hyperparameter search use a Bayesian regression model to estimate both the expected value of validation set error for each hyperparameter and the uncertainty around this expectation.

The learning rate has a U-shaped curve for training error: Tuning the parameters other than the learning rate requires monitoring both training and test error to see if the model is overfitting or underfitting, then adjusting its capacity.

6.1 manual tuning

the goal of hyperparameter search is usually to find the lowest generalization error subject to some runtime and memory budget.

When there are three or fewer hyperparameters, the common practice is to perform grid search, followed by a U-shaped curve when plotted as a function of one of the hyperparameters.

If there is a gap between train and test more efficient than a grid search, have time to tune only one hyperparameter, tune the learning rate. The effective capacity of the model is then gathering more data is not correct for the optimization problem. feasible, the only other way to improve generalization error is to improve the learning algorithm itself.

6.2 random search

random search is exponentially more efficient than a grid search.

6.5 model-based hyperparameter optimization

Most model-based algorithms for hyperparameter search use a Bayesian regression model to estimate both the expected value of validation set error for each hyperparameter and the uncertainty around this expectation.

The learning rate has a U-shaped curve for training error: Tuning the parameters other than the learning rate requires monitoring both training and test error to see if the model is overfitting or underfitting, then adjusting its capacity.

5. It is recommended to experiment with training set sizes on a logarithmic scale.

the learning rate requires monitoring both training and test error to see if the model is overfitting or underfitting, then adjusting its capacity.

7.1 test error is the sum of the training and test errors.

(vii) test error provides guarantees about the quality of the optimization process.

a smaller model that requires less memory and runtime to store and evaluate.

(ii) model parallelism: multiple machines work together on a single datapoint, with each machine running a different part of the model.

can be run by a separate machine. 2. implementations.

2.6 one strategy for accelerating a single datapoint, with data processing systems in general is to build systems that have dynamic structure

(iii) data parallelism is feasible in the graph describing the computation needed to process an input.

One major obstacle to using dynamically structured systems is the decreased degree of parallelism that results from the system following different code branches for different inputs.

(iv) data parallelism for training is somewhat harder, and can be solved using asynchronous stochastic gradient descent.

2.5 a key strategy for reducing the cost of inference is model compression.

The basic idea of model compression is to replace the original expensive model with

multiple CPU implementations. Careful implementation for specific CPU families can yield large improvements in runtime.

2.3 GPU implementations

GPUs are specialized hardware components that were originally developed for graphics applications. Both inference and training. General purpose GPUs could execute arbitrary code, not just performance GPU code, researchers should structure their workflow to avoid needing to write new and code branches for different inputs.

2.4 large-scale distributed implementations.

(i) data parallelism: each input example we want to process

Applications

Chapter 12

1. large-scale deep learning
 - ① deep learning is based on the philosophy of connectionism: while an individual biological neuron or an individual feature in a machine learning model is not intelligent, a large population of these neurons or features acting together can exhibit intelligent behavior.
 - ② One of the key factors responsible for the improvement in neural networks accuracy and the improvement of the complexity of networks they can solve between the 1980s and today is the dramatic increase in the size of the networks we use.
 - ③ Because the size of neural networks is of paramount importance, deep learning requires high performance hardware and software infrastructure.
- 2.7 specialized hardware implementations.

(v). Global contrast normalization (ACN) can be understood as mapping examples to a spherical shell, reducing each example to a direction rather than a direction and a distance. This is a useful property because neural networks are often better at responding to directions in space rather than exact locations.

(vi) Local contrast normalization ensures that the contrast is normalized across each small window, rather than over the image as a whole. It is a differentiable operation and can also be used as a nonlinearity applied to the hidden layers, as well as a preprocessing operation applied to the input.

is the only kind of pre-processing that is strictly necessary

different forms of specialized hardware have been developed over the last decades.

(ii) Rescaling (resizing) images Using fixed point rather than floating point representations is not always strictly necessary. Some convolutional and using less bits per number of the model to determine the output. This idea can be interpreted as an ensemble approach, and helps to reduce generalization error.

models accept variably-sized inputs and dynamically adjust the size of their pooling regions to keep the output size constant. Other convolutional models 3. computer vision.

have variable-sized output that automatically scales in size with the input.

(iii) dataset augmentation may be seen as a way of preprocessing the training set only. It is an excellent way to reduce

(i) formatting images to have the same scale, i.e., standardizing images so that their pixels all lie in the same range, of most computer vision models.

is the only kind of pre-processing that is strictly necessary

different hardware have been developed over the last decades.

(ii) Rescaling (resizing) images Using fixed point rather than floating point representations is not always strictly necessary. Some convolutional and using less bits per number of the model to determine the output. This idea can be interpreted as an ensemble approach, and helps to reduce generalization error.

models accept variably-sized inputs and dynamically adjust the size of their pooling regions to keep the output size constant. Other convolutional models 3. computer vision.

have variable-sized output that automatically scales in size with the input.

(iii) dataset augmentation may be seen as a way of preprocessing the training set only. It is an excellent way to reduce

(i) formatting images to have the same scale, i.e., standardizing images so that their pixels all lie in the same range, of most computer vision models.

between words that are in the same category.

5.4 Neural language models (NLNs) are a class of language models designed to overcome the curse of dimensionality problem for modeling natural language sequences by using a distributed representation of words.

Distributed representation can counter the curse of dimensionality.

5.5 high-dimensional outputs

- (i) using a softmax over a large number of output words has high cost. The first neural language model dealt with the high cost by using a short list. The obvious disadvantage of the short list approach is that the potential generalization of the neural language models is limited to the most frequent words, where, arguably, it is the least useful.
- (ii) To improve the statistical efficiency of n-gram models, a short list approach is that the class-based language models introduce the notion of word categories and then share statistical strength

the corresponding sequence of the nth token given the speaker.

Models based on n-grams define the conditional probability of words intended by the speaker.

The n-gram model uses products of these conditional distributions to define the probability distribution over longer sequences:

Distributed representation can counter the curse of dimensionality.

$$P(x_1, \dots, x_T) = P(x_1, \dots, x_N) \prod_{t=1}^T P(x_t | x_{t-1}, \dots, x_1)$$

$$= P(x_1, \dots, x_N) \prod_{t=1}^T P(x_t | x_{t-1}, \dots, x_1)$$

(ii) classical n-gram models are particularly vulnerable to the curse of dimensionality. To improve the statistical efficiency of n-gram models, a short list approach is that the class-based language models

introduce the notion of word categories and then share statistical strength

4. speech recognition
4.1 speech recognition

4.2 ASR: automatic speech recognition.

$$f_{ASR}^*(X) = \underset{Y}{\operatorname{argmax}} P(Y|X=X)$$

4.3 For AMT-HMM systems, AMMs modeled the association between acoustic features and phonemes, while HMMs are based on language models, modeled the sequence of phonemes.

4.4. GMM-HMM (1980s - 2012)
(i) A language model defines a probability distribution over sequences of tokens in a natural language. Depending on how the model is designed, some deep learning systems learn feature from raw acoustic input.

4.5 The task of speech recognition is to map an acoustic signal containing a spoken natural language utterance into

(i) An n-gram is a sequence of n tokens.

(ii) An n-gram is a spoken natural language utterance into

6. other applications

6.2 recommender systems.

Both online advertising and item recommendation systems :

(i) some components of machine rather than proportional to the number of negative samples (ii) the hierarchical softmax is introduced to the context of neural language model. It brings computational benefits both at training time and at test time, if at test time we want to compute the probability of many candidate translations.

(b). a language model which evaluates the proposed translations.

6.3 (i) Many recommendation problems are most accurately described theoretically as contextual bandits.

(iii) The encoder-decoder reinforcement learning can involve a sequence of many actions and many rewards. The bandits scenario is a special case of reinforcement learning, in which the learner takes only a single action and receives a single reward. The contextual bandits refers to the case where an action is taken in the context of some input variable that can inform the decision. The mapping from context to action is called a policy.

6.4. Knowledge of relations combined with a reasoning process and understanding of natural language could allow us to build a general question answering system.

meaning in another language. (i) the number of negative samples (ii) The hierarchical softmax size of the output vector.

5.6 combining neural language models with n-grams.

A major advantage of n-gram to compute the probability of models over neural networks specific words.

5.7. Neural machine translation

(i) machine translation is the task of reading a sentence in one language and emitting a weighted average. It is essentially a weighted average.

3.4 Linear factor models
including PCA and factor analysis can be interpreted as learning a manifold.

4. ICA

2.2 $x \sim N(x; b, W\tilde{W}^T + \psi)$
the role of the latent variables is to capture the dependencies between the different observed variables x_1 .

4.2 Independent component analysis (ICA)

is an approach to modeling linear factors that seeks to separate an observed signal into

many underlying signals that are scaled and added together

$x \sim N(x; b, W\tilde{W}^T + \sigma^2 I)$

to form the observed data. These signals are intended to be fully independent.

4.3 All variants of ICA require that $p(h)$ be non-Gaussian.

A generative model

either represents $p(x)$ or can draw samples from it.

some of the simplest generative models and some of the simplest probability models with latent variables

that learn a representation of data.

Much as linear classifiers and linear regression models may be extended to MRFs, these linear factor models may be extended to autoencoder networks and deep probabilistic models that perform the same tasks but with a much more powerful and flexible model family.

3.3 Probabilistic PCA model takes advantage of the observation that most

variations in the data

can be captured by the

latent variables h , up to

some small residual

reconstruction error σ^2 .

It becomes PCA as $\sigma \rightarrow 0$.

④ linear factor models are

some of the simplest generative models and some of the simplest probability models with latent variables

that learn a representation of data.

② A linear factor model is defined by

the use of a stochastic, linear decoder function that generates x by adding noise to a linear transformation

of h . Probabilistic PCA, factor

analysis, and many linear factor models that perform the same

tasks but with a much more

powerful and flexible

model family.

2. Factor analysis.

2.1 $h \sim N(h; 0, I)$

$x = Wh + b + \text{noise}$,

noise $\sim N(\text{noise}; 0, \phi)$

where $\phi = \text{diag}(\sigma^2)$, with

$\sigma^2 = [\sigma_1^2, \dots, \sigma_n^2]^T$

Chapter 13. Linear Factor Models

Chapter 13.

1.0 linear factor models are the simplest probability models with latent variables

that learn a representation of data.

② A linear factor model is defined by

the use of a stochastic, linear decoder

function that generates x by adding noise to a linear transformation

of h . Probabilistic PCA, factor

analysis, and many linear factor models that perform the same

tasks but with a much more

powerful and flexible

model family.

First, we sample the factors h from a distribution $h \sim p(h)$, where $p(h)$ is a factorial distribution, with $p(h) = \prod p(h_i)$,

Next we sample the real-valued observable variables given the factors:

$x = Wh + b + \text{noise}$, where the noise is typically Gaussian and diagonal.

6.3 The encoder that we use with sparse coding is not a parametric encoder. optimization problem

$$\min_{\theta} E(f(x^{(t)}), -f(c^{(t)}))^2$$

Instead, the encoder is an optimization algorithm, that solves an optimization problem in which we seek the single most likely code value: $h^* = f(x) = \arg \max p(h|x)$.

6.4 advantages

- (i) It can in principle minimize the combination of reconstruction error and log prior better than any specific parametric encoder.
- (ii) there is no generalization error to the encoder.

6.5 disadvantages.

- (i) the non-parametric encoder requires greater time to compute h given x .
- (ii) it is not straight-forward to back propagate through to nonparametric encoder.

6.6 sparse coding, like other linear factor models, often produces poor samples. This happens even when the model is able to reconstruct the data well and provide useful features for a classifier.

$$6.2 \quad i' s p(x|h) = \mathcal{N}(x; \mathbf{W}h + \beta)$$

where β is precision.

4.5 Many variants of ICA are to learn invariant features. ICA's more often used as an analysis tool for separating characteristics of scenes - change very slowly compared to the individual measurements that make up a description of its density.

5.3 The idea of slowness principle is that the important principle to any differentiable function f to the nonlinear autoencoders, model trained with gradient descent. SFA is that it is possibly to theoretically predict which feature SFA will learn, even in the deep, nonlinear setting. 5.5 SFA is not quite a generative model perse, in the sense that it defines a linear map from input space and feature space but doesn't define a prior over feature space and thus does not impose a prior on input space.

5.4 SFA is motivated by the slowness principle. 5.5 SFA is not quite a generative model perse, in the sense that it defines a linear map from input space and feature space but doesn't define a prior over feature space and thus does not impose a prior on input space.

5.6 The SFA algorithm consists of defining $f(x; \theta)$ to be

5.7 SFA is a linear factor model that uses info from time signals

is great enough to learn at most 2.3. The learning process is described simply as minimizing identity function.

3.3 nearly any generative model with latent variables and equipped with an inference procedure (for compacting latent representations given input) may be viewed as a particular form of autoencoder.

3.4 examples:

sparse autoencoder,
DAE,
CAE.

4. Sparse Autoencoders

A sparse autoencoder is simply an autoencoder whose training criterion includes a sparsity penalty $\Omega(h)$ on the code layer h , in addition to the reconstruction error: to become too great.

$$L(x, g(x)) + \Omega(h)$$

3. Regularized Autoencoders

4.3 We can think of the penalty $\Omega(h)$ as a regularizer term added to a feedforward network whose primary task even if the model capacity

Utilize general feedforward networks, autoencoders may also be trained using recirculation.

1.4 recirculation is a learning algorithm based on comparing the activations of the network on the original input to the activations on the reconstructed input. It is regarded as more biologically plausible than back-propagation, but is rarely used for machine learning applications.

2.5 An autoencoder can fail to learn anything useful about the dataset if the capacity of the autoencoder is allowed to become too great.

3. Undercomplete Autoencoders

3.2 A regularized autoencoder can be nonlinear and overcomplete to capture the most salient features of the training data.

Chapter 14 Autoencoders

1. An autoencoder is a neural network that is trained to attempt to copy its input to its output. Internally, it has a hidden layer h that decodes a code used to represent the input. The autoencoder has two components: an encoder function $f(x)$ and a decoder function $r = g(h)$.

1.2 Modern autoencoders have generalized the idea of an encoder and a decoder beyond deterministic functions to stochastic mappings $P_{\text{Encoder}}(h|x)$ and $P_{\text{Decoder}}(x|h)$.

1.3 Autoencoders may be thought of as being a special case of feedforward networks, and may be trained with all of the same techniques.

(ii) sample \tilde{x} from corruption

$$C(\tilde{x} | x = x)$$

(iii) use (\tilde{x}, x) as a training example for estimating the reconstruction distribution

$$\text{Procrustes}(x|\tilde{x})$$

$$= P_{\text{decoder}}(x|h), \text{ with } h(x)$$

(iii) This view provides a different motivation for training an autoencoder: it is just a way of approximately training for h . With this chosen h , we are maximizing

$$\log P_{\text{Model}}(h, x) = \log P_{\text{Model}}(h) + \log P_{\text{Model}}(x|h).$$

It also provides a different reason for why the features learned by autoencoder are

useful: they describe the latent variables that explain the input.

We can therefore view the DAE as performing stochastic gradient descent on the expectation:

5.2 A denoising autoencoder minimizes $L(x, g(f(x)))$, where \tilde{x} is a copy of x that has been corrupted by some form of noise.

5.5 (i) ~~Score~~ matching is alternative to maximum likelihood. It provides a consistent estimator of probability distributions based on encouraging

We can think of the autoencoder's to copy the input to the output. i: the regularizer depends on the data, i: there is not a straightforward Bayesian interpretation to this regularizer. But we can still think of these regularization terms as implicitly expressing a preference over functions.

The $\log P_{\text{Model}}(h)$ term can be sparsity-inducing.

4.4 (i) Rather than thinking of the sparsity penalty as a regularizer for the copying task, we can think of the entire sparse autoencoder framework as approximating maximum likelihood approximating maximum likelihood training of a generative model that has latent variables. (ii) For a model with visible variables x and latent variables h , with an explicit $P_{\text{Model}}(x, h) = P_{\text{Model}}(h)p_{\text{Model}}(x|h)$ the log-likelihood can be decomposed as

$$\log P_{\text{Model}}(x) = \log \sum_h P_{\text{Model}}(h, x).$$

9. stochastic encoders and decoders.

9.2 To make a more radical departure from feedforward networks, we can generalize the notion of an encoding function $f(x)$ to

an encoding distribution $P_{\text{encoder}}(h|x)$.

9.3 Any latent variable model $P_{\text{model}}(h|x)$ defines a stochastic encoder $P_{\text{encoder}}(h|x) = P_{\text{model}}(h|x)$

and a stochastic decoder $P_{\text{decoder}}(x|h) = P_{\text{model}}(x|h)$.

9.4 So long as the encoder is deterministic, the autocoder representing some functions is a feedforward network.

10. Learning manifolds with autocoders.

$$\min L(x, g(f(x))) + \Omega(h, x)$$

with $\Omega(h, x) = \lambda_2 \|\nabla_h h\|^2$

An autocoder regularized in this way is called contractive autocoder algorithm is equivalent to CCAE).

(RBMs with Gaussian visible units. When the RBM is trained using denoising training point x . In this context, score matching, its learning the score is a particular gradient field: $\nabla_x \log p(x)$)

(i) Learning the gradient field corresponding autocoder.

(ii) Learning the structure of Data itself.

(iii) A very important property of DAEs is that their training criterion (with conditionally Gaussian $p(x|h)$) makes the autoencoder learn a sparse autoencoders, sparse coding, contractive autoencoders and other regularized autoencoders,

5.6 like sparse autoencoders, encoder is itself a feedforward network as is the decoder. Advantages of depth in feedforward the motivation for DAEs networks

(with conditionally Gaussian $p(x|h)$)

6. Regularizing by penalizing derivatives.

(iv) Denoising training of a specific kind of autoencoder (sigmoidal hidden units, linear reconstruction units) using Gaussian noise and mse as the reconstruction cost is equivalent to training

13.

13.2 lower-dimensional representations can improve performance on many tasks.

Models of smaller spaces consume less memory and runtime.

12. PSD.

12.2 Predictive sparse

retrieval is to find entities in a database that resemble a query entry. This task derives the benefit that search can become extremely efficient in certain kinds of low dimensional spaces.

13.4 semantic hashing is an approach to information retrieval via dimensionality reduction and binarization.

11.7. The contraction penalty in the limit of small Gaussian noise can obtain useless results if we do not pose some sort of scale on the decoder.

Models of smaller spaces consume less memory and runtime.

12. PSD.

12.2 Predictive sparse

decomposition (PSD) is a model that is a hybrid of sparse coding and parametric autoencoders. Here $f(x)$ and $g(h)$ are both parametric. A parametric encoder is 11.5 CAE is contractive only trained to predict the output locally : all perturbations of iterative inference. a training point x are mapped near to $f(x)$.

12.3 training proceeds by minimizing $\|x - g(f(x))\|^2 + \gamma \|h - f(x)\|^2$. The goal of CAE is to learn the manifold structure of the 11. Contractive Autoencoders. 12.4 PSD is an example of learned approximate inference. (Chapter 19).

11.3. In the limit of small Gaussian noise, the denoising reconstruction error is equivalent to a contractive penalty on the reconstruction function that maps x to $z = g(f(x))$.

11.4. The contractive autoencoder introduces an explicit regularizer on the code $h = f(x) :$ $\mathcal{L}(h) = \lambda \left\| \frac{\partial f(x)}{\partial x} \right\|_F^2$ $\|\cdot\|_F$ is Frobenius norm.

10.3. The autoencoder can afford to represent only the variations that are needed to reconstruct training examples. The encoder learns a mapping that is only sensitive to changes along the manifold directions, but is insensitive to changes orthogonal to the manifold.

11.6 . The goal of CAE is to learn the manifold structure of the 11. Contractive Autoencoders.

2.4 unsupervised pretraining combines two different ideas.

First, it makes use of the idea that the choice of initial parameters for a dm can have a significant regularizing effect on the model.

Second, it makes use of the more general idea that learning about the input distribution can help to learn

about the mapping from inputs to outputs.

2.5 The first idea is least well understood. Our ability to characterize exactly what aspects of the pretrained could be used to find a good initialization for a joint learning parameters are retained during the supervised training stage is limited.

This is one reason that modern approaches typically use simultaneous unsupervised learning even fully connected and supervised learning rather than two sequential stages.

2.6. The other idea is better understood. The basic idea is that some features that are

2.1 unsupervised pretraining, or

more precisely, greedy unsupervised pretraining, played a historical role

in the revival of dm, enabling

researchers for the first time to train a deep supervised network without requiring architectural specializations like convolution or recurrence.

2.2 The deep learning renaissance

of 2006 began with the discovery that this greedy learning procedure could explain improved human performance—for example, the

brain may use very large ensembles of classifiers or Bayesian inference techniques.

One popular hypothesis is that the brain is able to leverage unsupervised or semi-supervised learning.

2.3 unsupervised pretraining is sometimes helpful but often harmful.

1.4 most representation learning problems face a tradeoff between

preserving as much information about the input as possible and obtaining nice properties.

1.5 We can learn good representations for the unlabeled data, and then use these representations to

solve the supervised learning task.

1.6. Humans and animals are able to learn from very few labeled examples. We do not yet know how this is possible.

Many factors representation at every hidden layer taking on properties that make the classification task easier.

1.3 supervised training of feedforward networks does not involve explicitly imposing any condition on the learned intermediate features. Other kinds of representation learning algorithm

shape the representation in some particular way.

2. greedy layer-wise unsupervised pretraining

Chapter 15 Representation learning

1. Many information processing tasks can be very easy or very difficult depending on how the info is

represented. A good representation

is one that makes a subsequent learning task easier.

1.2 We can think of feedforward networks trained by supervised learning as

performing a kind of representation learning. Training with a supervised criterion naturally leads to the

best possible. Many factors representation at every hidden layer

taking on properties that make the classification task easier.

1.3 supervised training of feedforward networks does not involve explicitly

imposing any condition on the learned intermediate features. Other kinds of representation learning algorithm

shape the representation in some particular way.

2. greedy layer-wise unsupervised pretraining

3.5 concept drift can be viewed as a form of transfer learning, due to gradual changes in the data distribution over time.

3.6 Both concept drift and transfer learning situation where what has been learned in one setting (i.e., distribution P_1) is exploited to improve generalization in another setting (say distribution P_2).

3.7 The core idea of representation learning is that the same representation may be useful in both settings.

3.8 two extreme forms of transfer learning are one-shot learning and zero-shot learning, sometimes also called zero-data learning.

4. semi-supervised disentangling of causal factors.

4.2. an ideal representation is one that disentangles the causes from one another.

3. Transfer learning and domain adaptation.

3.2 transfer learning and domain adaptation refer to the situation where what has been learned in one setting (i.e., distribution P_1) is exploited to

another setting (say distribution P_2).

3.3 In transfer learning, the learner must perform two or more different tasks, but we assume that many of the features with dropout or batch normalization, that explain the variations in P_1 , are relevant to the variations that need to be captured for learning P_2 .

3.4 In the related case of domain adaptation, the task remains the same between each setting, but the input distribution is slightly different.

2.9. disadvantages of unsupervised pretraining.

(i) there is not a way of flexibly adapting the strength of the regularization arising from the unsupervised stage.

(ii) Another disadvantage of having two separate training phases is that each phase has its own hyperparameters.

2.10. deep learning techniques based on supervised learning, regularized on simultaneous supervised and unsupervised learning can be prefferable.

2.7 unsupervised pretraining can act on extremely small datasets, Bayesian methods outperform methods based on unsupervised pretraining.

2.8 how can unsupervised pretraining act as a regularizer?
One hypothesis is that pretraining encourages the learning algorithm to discover features that relate to the underlying causes that generate the observed data.

2.11 The idea of pretraining has been generalised to supervised pretraining, as a very common approach for transfer learning.

one is that generalization arises due to shared attributes between different concepts.

5.6 **Distributed representations introduce a rich similarity space**

5.7. When and why can there be a statistical advantage from using a distributed representation as part of a learning algorithm?

i) Distributed representations can have a statistical advantage when an apparently complicated structure can be compactly represented using a small number of parameters.

(ii) with $O(nd)$ parameters (for n linear-threshold features in R^d) we can distinctly represent $O(nd)$ regions in input space. If instead we used a kernel, made no assumption at all about the data, and used a representation based on n-grams.

to describe k^n different concepts

5.3 A symbolic representation is a specific example of the broader class of non-distributed

representations, which are representations that may contain many entries but without significant meaningful separate control over each entry.

5.4 examples of learning algorithms based on non-distributed representations:

clustering methods,
k-nearest neighbors algorithms,

decision trees,

Gaussian mixtures and mixtures of experts (in $P(y)$).

5. Distributed representation .

5.1 distributed representations composed of many elements that can be set separately from each other. They are powerful because they can use n features with k values

We expect that future research will develop mechanisms for representing different factors depending on the task.

4.3 A representation that clearly separates the underlying causal factors may not necessarily be one that is easy to model. If a representation h represents many of the underlying causes of the observed x , and the outputs y are among the most salient causes, then it is easy to predict y from h .

4.4. human beings fail to perceive changes in their environment that are not immediately relevant to the task they are performing.

4.5 An important research frontier in semi-supervised learning is determining what to encode in each situation. GRNs are only one step toward determining which factors should be represented.

learning provides a very strong
due to: a label that specifies
the value of at least one
of the factors of variation
directly. Representation learning
~~uses~~ makes use of other,
less direct hints about the
underlying factors.

models with a single hidden
layer of latent variables,
including RBMs and deep
belief networks, are
universal approximators
of probabilistic
probability distributions

7.2 An ideal representation is
one that disentangles the
underlying causal factors of
variation that generated the
data, especially those factors
that are relevant to our
applications. Most strategies
for representation learning
are based on introducing class
factors, shared factors across tasks,
manifolds, natural clustering,

such as non-distributed
models generalize only locally
via the smoothness prior. (iii) A further part of the argument
of the factors of variation
directly. Representation learning
~~uses~~ makes use of other,
less direct hints about the
underlying factors.

6.2 It has been proven in many
different setting that organizing
computation through the
composition of many
nonlinearities and a
hierarchy of reused features
can give an exponential
boost to statistical efficiency!
on top of the exponential
boost given by using a
distributed representations.

6.3 deterministic feedforward
networks are universal
approximators of functions.
Many structured probabilistic

then specifying $O(n^d)$ regions
would require $O(n^d)$ examples.
representations generalize well if
they encode so many different regions.

5.8 Non-distributed algorithms break the
input space into regions, with a
separate set of parameters for
each region. The advantage of
a non-distributed approach is that,
given enough parameters, it can
fit the training set without
solving a difficult optimization
algorithm, because it is
straightforward to choose a different
output independently for each
region. The disadvantage is that

learning provides a very strong
due to: a label that specifies
the value of at least one
of the factors of variation
directly. Representation learning
~~uses~~ makes use of other,
less direct hints about the
underlying factors.

6. Exponential gains from depth
representations generalize well if
that their capacity remains limited
despite being able to distinctly

such as non-distributed
models generalize only locally
via the smoothness prior. (iii) A further part of the argument
of the factors of variation
directly. Representation learning
~~uses~~ makes use of other,
less direct hints about the
underlying factors.

4.4 The graph encodes only simplifying assumptions about which variables are conditionally independent from each other.

It is also possible to make other kinds of simplifying assumptions.

4.5 The directed graphical model syntax does not place any constraint on how we define our conditional distributions. It only defines which variables they are allowed to take in as arguments.

5. undirected models

5.2 undirected models also known as Markov random fields (MRFs) or Markov networks.

5.3 unnormalized probability distribution.

$$\tilde{p}(x) = \prod_{C \in G} \phi(C)$$

where C is clique in the graph, where $\tilde{p}_G(x)$ gives the factor $\phi(C)$ is clique potential.

3. Using graphs to describe model structure

3.1 structured probabilistic models use graphs to represent then the naive approach of representing $p(x)$ by storing a lookup table with one probability per possible outcome

Each edge represents a direct interaction.

3.2 graphical models can be largely divided into two categories: models based on directed acyclic graphs, and models based on undirected graphs.

4. directed models

4.2 directed graphical models, also known as belief network or Bayesian network.

$$p(x) = \prod p(x_i | \text{Par}(x_i))$$

where $\text{Par}(x)$ gives the parents of x in G .

2.5 modeling model a distribution over a random vector x containing n discrete variables capable of taking on k values each,

representing $p(x)$ by storing a lookup table with one probability per possible outcome requires k^n parameters.

This is not feasible for several reasons:

- (i) memory: the cost of storing the representation.
- (ii) statistical efficiency.
- (iii) runtime: the cost of inference structure.

2.6 The table-based approach

2.3 classification algorithms are often able to ignore many parts of the input.

2.4 Many other tasks are often more expensive than classification. Most require a complete understanding of the entire structure of the input, with no option to ignore sections of it. These tasks include the following:

density estimation, denoising, missing value imputation, sampling.

Chapter 16 Graphical models for AI

(Graphical model) is a way of describing a probability distribution, using a graph to describe which random variables in the probability distribution interact with each other directly.

2. The challenge of unstructured modeling

- 2.1 The goal of deep learning is to scale machine learning to the kinds of challenges needed to solve AI. This means being able to understand high-dimensional data with rich structure.

2.2 The challenge of unstructured modeling

2.3 classification algorithms are often able to ignore many parts of the input.

2.4 Many other tasks are often more expensive than classification. Most require a complete understanding of the entire structure of the input, with no option to ignore sections of it. These tasks include the following:

density estimation, denoising, missing value imputation, sampling.

8.3 two variables are dependent just a special kind of Markov network : the exponentiation makes each term in the energy function correspond to a factor for a different clique.

8.4. In the case of undirected models, we refer to paths involving only unobserved variables as "active" and paths including an observed variable as "inactive".

$$F(x) = -\log \sum_h \exp(-E(x, h))$$

8.5 In directed models, all kinds of active paths of length two that can exist between a and b .

- (i) $\textcircled{a} \rightarrow \textcircled{b} \rightarrow \textcircled{d}$
- (ii) $\textcircled{a} \leftarrow \textcircled{b} \leftarrow \textcircled{d}$
- (iii) $\textcircled{a} \rightarrow \textcircled{b} \rightarrow \textcircled{d}$
- (iv) $\textcircled{a} \rightarrow \textcircled{b} \leftarrow \textcircled{d}$

here the shaded variables are "dependent," indicating that they are observed.

7. Energy-based models (EBMs)
- 7.1 An energy-based model is just a special kind of Markov network : the exponentiation makes each term in the energy function correspond to a factor for a different clique.
- 7.2 Many interesting theoretical results about undirected models depend on the assumption that $\tilde{P}(x) > 0$, i.e. A convenient way to enforce this condition is to use an EBM where
- $$\tilde{P}(x) = \exp(-E(x)), \quad (1)$$
- and $E(x)$ is known as the energy function.
- 7.3 Any distribution of the form given by equation (1) is an example of a Boltzmann distribution. For this reason, many energy-based models are called Boltzmann machines.
- 7.4. The term Boltzmann machine probability distributions.
- 7.5 It is possible to leverage the effect of a carefully chosen domain of a variable in order to obtain complicated behavior from a relatively simple set of ϕ functions.

generally very fast (assuming sampling from each conditional is easy) and convenient. The drawbacks are that it only applies to directed models and does not support every conditional sampling operation.

11.4 Sampling techniques for undirected models are an advanced topic, covered in more detail in Chapter 17.

12. Advantages of structured modeling

12.2 (i) dramatically reduce the cost of representing probability distributions as well as learning and inference.

(ii) sampling is accelerated in the case of directed models.
(iii) allow us to explicitly separate representation of knowledge from learning or knowledge or inference given existing knowledge. This makes our models easier to develop and debug.

10. Factor graphs

10.2 Factor graphs are another way of drawing undirected models that resolve an ambiguity in the graphical representation of standard undirected model syntax.

10.3 factor graphs explicitly represent the scope of each function.

we can always represent any distribution by using a complete graph.

9.4 When we represent a probability distribution with a graph, we want to choose a graph that implies as many independencies as possible without implying any independencies that do not actually exist.

9.5 Directed models are able to use one specific kind of substructure that undirected models cannot represent perfectly. This substructure is called an immorality.

9.6 Likewise, undirected models described using an undirected no directed model can represent perfectly.

11.3 ancestral sampling is

8.6 context-specific independencies are independencies that are present dependent on the value of some variables in the network.

8.7 In general, a graph will never imply that an independence exists when it does not. However, a graph may fail to encode an independence.

9. converting between undirected and directed graphs.

9.2 no probabilistic model is inherently directed or undirected. Instead, some models are most easily described using a directed graph, or most easily described using an undirected graph.

9.3 Every probability distribution can be represented by either a directed model or by an undirected model. In the worst case,

- 15.6. Most deep models are designed to make Gibbs sampling or variational inference algorithms efficient.
- 15.7. the deep learning approach to graphical modeling is characterized by a marked tolerance of the unknown, always makes use of the idea of distributed representations.
16. The restricted Boltzmann machine (RBM)
- 16.2. The canonical RBM is an energy-based model with binary visible and hidden units. Its energy function is
- $$E(v, h) = -\sum_v v^T \phi_v - \sum_h h^T \phi_h$$
-
- 16.3. RBM demonstrates the typical deep learning approach to graphical models: representation learning can be completed between two layers of latent variables, combined via layers of interactions between groups may be described with efficient interactions between layers parametrized by matrices.

in order to implement a learning rule.

we can define the depth of a model in terms of the graphical model rather than the computational graph. 15.3. Deep learning essentially learning are usually not restrictive enough to allow efficient inference. This motivates the use of approximate inference. (See Chapter 19).

13. Learning about dependencies of inference problems in which we must predict the value of some variables given other variables, or predict the probability distribution over some variables given the value of other variables.
- 14.3. The graphs used for deep learning are usually not restrictive enough to allow efficient inference. This motivates the use of approximate inference. (See Chapter 19).
15. The deep learning approach to graphical models typically have large groups of units that are connected to other groups of units, so that representation learning can be completed between two layers of latent variables, combined via layers of interactions between groups may be described with efficient interactions between layers parametrized by matrices.

13. Learning about dependencies of inference problems in which we must predict the value of other. In the context of deep learning, hidden variables have been introduced to model these dependencies. 13.3. Latent variables have advantages beyond their role in efficiently capturing $p(v)$. The new variables h also provide an alternative to allow efficient inference. 14.2. Because $\log p(v) = E_{h \sim p(h)} (\log p(h, v) - \log p(h|v))$, we often want to compute $p(h|v)$,

3.5 biased importance sampling to confront the intractable partition function of untrained models. (See Chapter 8).

$$\hat{S}_{\text{BIAS}} = \frac{\sum_{i=1}^n \frac{p(x^{(i)})}{q(x^{(i)})} f(x^{(i)})}{\sum_{i=1}^n \frac{f(x^{(i)})}{q(x^{(i)})}}$$

where \hat{p} and \hat{q} are the unnormalized forms of p and q and the $x^{(i)}$ are the samples from q . This estimator is biased because $E[\hat{S}_{\text{BIAS}}] \neq S$, but is asymptotically unbiased.

4. MCMC methods

4.2 we have use Markov chain to approximately sample from $P_{\text{Model}}(x)$. The family of algs that use Markov chains to perform Monte Carlo estimates

$$\hat{S}_Q = \frac{1}{n} \sum_{i=1}^n \frac{p(x^{(i)})}{q(x^{(i)})} f(x^{(i)})$$

is called Markov chain Monte Carlo methods (MCMC).

4.3 The most standard, generic guarantees for MCMC techniques are only applicable when the model does not assign zero probability to any state.

to confront the intractable partition function of untrained models. (See Chapter 8).

3. importance sampling

3.2 Importance sampling and its variants have been found very useful in many machine learning algs, including all alg.

3.3 `` $S = \int p(x) f(x) dx = E_p[f(x)]$ where $x^{(i)}$ are samples from p .
 $f(x) = q(x) \frac{p(x)}{q(x)}$

any Monte Carlo estimator

$\hat{S}_p = \frac{1}{n} \sum_{i=1}^n \frac{p(x^{(i)})}{q(x^{(i)})} f(x^{(i)})$

3.4 ① $E[\hat{S}_p] = E_p[\hat{S}_p] = S$

② $\text{Var}[\hat{S}_p] = \text{Var}\left[\frac{p(x^{(i)})}{q(x^{(i)})} f(x^{(i)})\right] / n$

3.5 minimum variance occurs where q is $q^*(x) = \frac{p(x)f(x)}{Z}$, where Z is the normalization constant.

Better importance sampling distribution
 put more weight where the integrand is larger.

view a sum or integral as if

it was an expectation under some distribution and to approximate the expectation by a corresponding average.

3.1 randomization fail into two rough categories: Las Vegas algorithms and Monte Carlo algorithms.

1.3 Las Vegas algorithms always return precisely the correct answer (or report that they failed).

1.4. Monte Carlo algorithms return answers with a random amount of error. The amount of error can typically be reduced by expending more resources.

1.5 many problems in machine learning are so difficult that we can never expect to obtain precise answers to them. This excludes precise deterministic alg and Las Vegas alg. Instead, we must use deterministic approximate alg or Monte Carlo approximations. Both approaches are ubiquitous in machine learning.

2.6 properties

① $E[\hat{S}_n] = S$

② if $x^{(i)}$ are i.i.d., $\text{Var}[f(x^{(i)})]$ is bounded, then

(i) $\text{Var}[\hat{S}_n] = \frac{1}{n} \text{Var}[f(x^{(i)})]$

= $\frac{\text{Var}[f(x)]}{n}$

(ii) $\lim_{n \rightarrow \infty} \hat{S}_n = S$, almost surely.

(iii) the distribution of \hat{S}_n converges to a normal distribution

with mean S and variance $\frac{\text{Var}[f(x)]}{n}$

2. sampling and Monte Carlo methods

2.3 why sampling?

It provides a flexible way to approximate many sums and integrals at reduced cost. In many cases, it is actually our goal, in the sense that we want to train a model that can sample from the training distribution.

2.4. The idea of Monte Carlo sampling is to

Chapter 7 Monte Carlo methods

1. randomized algorithms fail into two rough categories: Las Vegas algorithms and Monte Carlo algorithms.

1.3 Las Vegas algorithms always return precisely the correct answer (or report that they failed).

1.4. Monte Carlo algorithms return answers with a random amount of error. The amount of error can typically be reduced by expending more resources.

1.5 many problems in machine learning are so difficult that we can never expect to obtain precise answers to them. This excludes precise deterministic alg and Las Vegas alg. Instead, we must use deterministic approximate alg or Monte Carlo approximations. Both approaches are ubiquitous in machine learning.

2.6 properties

① $E[\hat{S}_n] = S$

② if $x^{(i)}$ are i.i.d., $\text{Var}[f(x^{(i)})]$ is bounded, then

(i) $\text{Var}[\hat{S}_n] = \frac{1}{n} \text{Var}[f(x^{(i)})]$

= $\frac{\text{Var}[f(x)]}{n}$

(ii) $\lim_{n \rightarrow \infty} \hat{S}_n = S$, almost surely.

(iii) the distribution of \hat{S}_n converges to a normal distribution

with mean S and variance $\frac{\text{Var}[f(x)]}{n}$

2. sampling and Monte Carlo methods

2.3 why sampling?

It provides a flexible way to approximate many sums and integrals at reduced cost. In many cases, it is actually our goal, in the sense that we want to train a model that can sample from the training distribution.

5. Gibbs sampling

4.9 running the Markov chain We describe q using a vector v , with 4.11 The core idea of a Markov chain until it reaches its equilibrium $q(x=i) = v_i$, we define A so that is to have a state x that begins as an arbitrary value. Over time, we $A_{ij} = T(x'=j|x=j)$.
 5.2 Now we have described how distribution is called "burning to draw samples from a distribution in" the Markov chain.
 $q(x)$ by repeatedly updating

$$x \leftarrow x' \sim T(x'|x).$$

In order to ensure $q(x)$ is a useful distribution, there are two basic approaches. The first one is to derive T from a given learned Pmodel (example 5.3).

The second one is to directly parametrize T and learn it, so that its stationary distribution implicitly defines the Pmodel of interest.

5.3 A conceptually simple and effective approach to building a Markov chain that samples from an energy-based model defining a distribution $P_{\text{model}}(x)$ is to use Gibbs sampling, in which sampling from $T(x'|x)$ is accomplished by selecting one

4.9 running the Markov chain until it reaches its equilibrium $q(x=i) = v_i$, we define A so that $A_{ij} = T(x'=j|x=j)$.
 5.2 Now we have described how distribution is called "burning to draw samples from a distribution in" the Markov chain.
 $q(x)$ by repeatedly updating

$$\therefore q^{(t+1)}(x') = \sum_j q^{(t)}(x) T(x'|x). \quad \text{Eventually, } x \text{ becomes (very nearly) a fair sample from } p(x).$$

4.10 Markov chains are expensive to use because of the time required to burn in the equilibrium distribution and the time required to transition from one sample to another reasonably decorrelated sample P we wish to sample from.

4.11 Suppose x has countably many states, then we can represent the states as just a positive integer x . We run infinitely many Markov chains in parallel. All of the states of the different Markov chains are drawn from some distribution $q^{(t)}(x)$, where t indicates the number of time steps that have elapsed. Our goal is for $q^{(t)}(x)$ to converge to $p(x)$.

4.12 Most properties of Markov chains can be generalized to continuous variables. In general, a Markov chain with transition operator T will converge, under mild conditions, to a fixed point described by the equation

$$q^*(x') = \text{Ex}_q T(x'|x)$$

- (i) $P_p(x) \propto \exp(-\beta E(x))$ it well, which implies that h they have a tendency to variables x_i and sampling it mix poorly.
- where P is described as being the mutual info. We often learn reciprocal of the temperature, reflecting the origin of energy-based generative models that very precisely samples become very correlated. We refer to such behavior as the energy-based model. It is also possible to sample several variables at the same time so long as they are conditionally independent given all of their neighbors. Gibbs sampling approaches.
- 6.3 In high dimensional cases, MCMC neighbors. Gibbs sampling approaches.
- models in the undirected graph G defining the structure of the energy-based model. It is also possible to sample several variables at the same time so long as they are conditionally independent given all of their neighbors. Gibbs sampling approaches.
- 6.4. In the context of models with latent variables, which define a joint distribution $P_{\text{model}}(x, h)$, we often draw samples of x by alternating between sampling from called block Gibbs sampling.
- 6.5 In many classification problems, we expect a distribution that has a manifold structure from high-temperature distributions with a separate manifold for each class: the distribution is concentrated around many modes, then resume sampling from the unit temperature distribution. These modes are separated by vast regions of high energy. This type of distribution would make MCMC methods converge very slowly because of poor mixing between modes.
- Another approach is to use parallel tempering.
- 6.6. Tempering to mix between representations that encodes x into h in such a way that a Markov chain in the space of h can mix more easily.
- 6.7 depth may help mixing One way to resolve the problem of 6.4 is to make h be a deep representation that encodes x into h in such a way that a Markov chain in the space of h can mix more easily.
- 6.8. The challenge of mixing between separated modes. The primary difficulty involved with MCMC methods is that

⑤ $g \leftarrow g - \frac{1}{m} \sum \nabla_{\theta} \log \tilde{p}(x^{(i)}; \theta)$
 partition function using gradient ascent:

⑥ $\theta \leftarrow \theta + \epsilon g$
 end while.

3.4. The contrastive divergence (CD) algorithm :

(i) Set ϵ , the step size, to a small positive number.

(ii) Set k , the number of Gibbs steps, high enough to allow burn in. Perhaps 100 to train an RBM on a small image patch.

(iii) while not converged do
 ① sample a minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from the training set.
 // ①, ②, ③, ④, ⑤, ⑥ are same as 3.3

// ③ : for $i=1$ to m do
 $x^{(i)} \leftarrow x^{(i)}$

end for

3.5 The CD algorithm initializes the Markov chain at each step with samples from the data distribution.

partition function using gradient ascent:

(i) Set ϵ , the step size, to a small positive number.

(ii) Set k , the number of Gibbs steps, high enough to allow burn in. Perhaps 100 to train an RBM on a small image patch.

(iii) while not converged do
 ① sample a minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from the training set.

② $g \leftarrow \frac{1}{m} \sum \nabla_{\theta} \log \tilde{p}(x^{(i)}; \theta)$
 ③ Initialize a set of m samples $\{x^{(1)}, \dots, x^{(m)}\}$ to random values.

④ for $i=1$ to k do
 for $j=1$ to m do
 $x^{(j)} \leftarrow \text{gibbs_update}(x^{(j)})$

end for
 end for
 end while

regularity conditions on \tilde{p} and $\nabla_{\theta} \tilde{p}(x)$,

$\nabla_{\theta} \log Z = \mathbb{E}_{x \sim p(x)} \nabla_{\theta} \log \tilde{p}(x)$. 1. We must normalize \tilde{p} by dividing by a partition function $Z(\theta)$ in order to obtain a valid probability distribution:

$$p(x; \theta) = \frac{1}{Z(\theta)} \tilde{p}(x; \theta)$$

This identity is the basis for a variety of Monte Carlo methods for approximately maximizing the likelihood of models with intractable partition functions.

2.4 The Monte-Carlo approach to learning undirected models provides an intuitive framework in which we can interpret

2. The log-likelihood gradient the positive phase as pushing down on the energy of training examples and the negative phase as pushing up the energy of samples drawn from the model.

2.2 The gradient of the log-likelihood of the partition function:

$\nabla_{\theta} \log p(x; \theta) = \nabla_{\theta} \log \tilde{p}(x; \theta) - \nabla_{\theta} \log Z$

3. CD and SMN 3.3 a naive MCMC algorithm for maximizing the log-likelihood with an intractable $p(x)$ for all x , and under

Chapter 18] confronting the partition function

regularity conditions on \tilde{p} and $\nabla_{\theta} \tilde{p}(x)$,

2.3 For models that guarantee

hidden and visible units.

$$4.2 \quad \frac{p(x)}{p(y)} = \frac{\sum \tilde{p}(x)}{\sum \tilde{p}(y)} = \frac{\tilde{p}(x)}{\tilde{p}(y)}$$

CD is only able to provide an initialization for the visible units. SML is able to train the pseudolikelihood is based on the observation that conditional probabilities take this ratio-based form, and thus can be computed without knowledge of the partition function.

4.3 $p(a|b) = \frac{p(a,b)}{p(b)} = \frac{\sum_{a,c} p(a,b,c)}{\sum_{a,c} p(a,b,c)}$ mix between steps.

$$\therefore p(x_i|x_{-i}) = \frac{\tilde{p}(x_i, x_{-i})}{\sum_{x_i} \tilde{p}(x_i, x_{-i})}$$

4.4. ~~likelihood~~ objective function:

$$\log p(x) = \log p(x_1) + \log p(x_2|x_1) + \dots + \log p(x_n|x_1, \dots, x_{n-1})$$

② pseudolikelihood objective function:

$$\sum \log p(x_i|x_{-i})$$

3.11. The SML/PCD algorithm using gradient ascent as the optimization procedure:

(1) Set ϵ

(ii) Set k , high enough to allow spurious modes.

a Markov chain sampling from 3.7. It was shown experimentally that the CD estimator is biased from samples from $p(x; \theta + \epsilon g)$ to burn in, starting that the CD estimator is biased for RBMs and fully visible Boltzmann machines in that it converges to a different point than the maximum likelihood estimator.

3.12 It is necessary to draw random values.

the samples starting from a fresh Markov chain initialized from a random starting point

// same as 3.3 after the model is done training. end while

① likelihood objective function:

$$\log p(x) = \log p(x_1) + \log p(x_2|x_1) + \dots + \log p(x_n|x_1, \dots, x_{n-1})$$

② pseudolikelihood objective function:

$$\sum \log p(x_i|x_{-i})$$

3.13 The regions that have high probability under the model but low probability under the data generating distribution are called spurious modes.

3.14 The regions that have high probability under the model but low probability under the data generating distribution are called spurious modes.

3.15 The regions that have high probability under the model but low probability under the data generating distribution are called spurious modes.

3.16 The regions that have high probability under the model but low probability under the data generating distribution are called spurious modes.

3.17 SML is considerably more resistant to forming models with spurious modes than contrastive divergence (PCD).

3.18 SML provides an initialization point for both the

where c is explicitly introduced as an approximation of $-\log Z(\theta)$.

7.3 NCE works by reducing training example.

the unsupervised learning problem of estimating $p(x)$ to that of learning a probabilistic binary classifier in which one of the categories corresponds to the data generated by the model.

7.4 we use standard maximum likelihood learning on the supervised learning problem of fitting Point to Point.

$$\theta, c = \underset{\theta, c}{\operatorname{argmax}} \sum_{x, y \in \text{train}} \log p_{\text{model}}(y|x)$$

7.2 In NCE, the probability where $P(y=1) = \frac{1}{2}$

7.3 minimizing the expected value of

$$P_{\text{Joint}}(x|y=1) = P_{\text{model}}(x)$$

$P_{\text{Joint}}(x|y=0) = P_{\text{model}}(x)$

$\log p_{\text{model}}(x) = \log \tilde{p}_{\text{model}}(x) + c$

$\log \tilde{p}_{\text{model}}(x) = P_{\text{Joint}}(x|y=0) - P_{\text{Joint}}(x|y=1)$

$\log \tilde{p}_{\text{model}}(x) = P_{\text{Joint}}(x|y=0) - P_{\text{Joint}}(x|y=1)$

$\log \tilde{p}_{\text{model}}(x) = P_{\text{Joint}}(x|y=0) - P_{\text{Joint}}(x|y=1)$

states that have only one variable different from a

training example.

6. denoising score matching

6.2 we fit a distribution

that corresponds to the true p_{data} ,

rather than the true p_{data} ,

where $q(x|y)$ is a corruption

to the data generated by the process.

6.3 several autoencoder

training algorithms are

equivalent to score matching

or denoising score matching.

7. Noise-contrastive estimation

(NCE)

7.2 In NCE, the probability

distribution estimated by the

model is represented explicitly

the bit at position j flipped,

minimizing the expected value of

5.4 score matching is not applicable to models of discrete data.

5.5. score matching can be viewed as a version of contrastive divergence using a specific kind of Markov

chain. The Markov chain in this

case is not Gibbs sampling, but

rather a different approach

that makes local moves

guided by the gradient.

5.6 ratio matching applies

specifically to binary data.

5.7 As with the pseudolikelihood

$L(x, \theta) = \sum_{j=1}^n \left(\frac{1}{1 + \frac{p_{\text{model}}(x_j, \theta)}{p_{\text{model}}(f(x_j), \theta)}} \right)^2$

5.8 score matching

$L(x, \theta) = \frac{1}{2} \| \nabla_x \log p_{\text{model}}(x; \theta) \|_2^2$

5.9 score matching

$L(x, \theta) = \frac{1}{2} \sum_j \left(\frac{\partial}{\partial x_j} \log p_{\text{model}}(x_j; \theta) \right)^2$

5.10 score matching

$L(x, \theta) = \frac{1}{2} \sum_j \left(\frac{\partial}{\partial x_j} \log p_{\text{model}}(x_j; \theta) \right)^2$

5.11 score matching

5.12 score matching

4.5 estimation by maximizing the pseudolikelihood is asymptotically consistent.

4.6 pseudolikelihood estimates can be thought of as having something resembling a negative phase.

5. score matching and ratio

matching

5.1 score matching

5.2 score matching

5.3 score matching

5.4 score matching

5.5 score matching

5.6 score matching

5.7 score matching

5.8 score matching

5.9 score matching

5.10 score matching

5.11 score matching

5.12 score matching

5.13 score matching

5.14 score matching

5.15 score matching

being too far from p_1 by

introducing intermediate distributions $= \sum_{j=0}^n p_j(x) \frac{p_i(x)}{p_j(x)} dx$
that attempt to bridge the gap
between p_0 and p_1 .

8.6 AJS:

$$\hat{\Sigma}_1 = \frac{Z_1}{Z_0} - \frac{Z_0}{Z_1} \dots - \frac{Z_{n-1}}{Z_n} \frac{Z_n}{Z_{n-1}} \dots \frac{Z_1}{Z_{n-1}}$$

s.t.: $x^{(k)} \sim p_0$
 $\frac{Z_1}{Z_0} = \frac{Z_1}{Z_1} \cdot \frac{Z_0}{Z_1} \dots \frac{Z_{n-1}}{Z_{n-2}} \frac{Z_n}{Z_{n-1}}$ 8.3 an estimate of the

variance of $\hat{\Sigma}_1$:

$$= \frac{n-1}{T} \frac{Z_{1:n}}{Z_1}$$

where a sequence of distributions

p_0, \dots, p_{n-1} , with $0 = n_0 < n_1 < \dots < n_{n-1}$ if p_0 is close to p_1 ,
8.2 can be an effective

way of estimating the
partition function. If not,

Provided the distributions p_j and
~~p_{j+1}~~ p_{j+1} , for $0 \leq j \leq n-1$, are
sufficiently close, we can
reliably estimate each of

the factors $\frac{Z_{j+1}}{Z_j}$ using simple 8.5 annealed importance sampling
importance sampling and then use (AIS) and bridge sampling
these to obtain an estimate to overcome the

problem of the proposal p
of $\frac{Z_1}{Z_0}$.

should be able to distinguish
data from noise. CRNs are
based on the idea that a
good generative model should

be able to generate samples
that no classifier can
distinguish from data.

7.5 Point is essentially a
logistic regression model.

7.8 SMLE and CP estimate
only the gradient of the
 $\hat{V}_{\text{ar}}(\hat{\Sigma}_1) = \frac{Z_0}{K^2} \frac{K}{\left(\frac{p_1(x^{(1)})}{p_0(x^{(1)})} - \frac{1}{2} \right)^2}$ log partition function. Score
matching and pseudolikelihood

$= \sigma((\log p_{\text{model}}(x) - \log p_{\text{noise}}(x))$
avoid computing quantities
related to the partition

function altogether.

7.6 NCE is most successful
when applied to problems with
few random variables, but

can work well even if those
random variables can take
on a high number of values

7.7 NCE is based on the idea
that a good generative model

3.6. Even though the E-step involves exact inference, we can think of the EM algorithm as using approximate inference in some sense.

3.7. insights.

First, there is the basic structure of the learning process, in which we update the model parameters to improve the likelihood of a completed dataset, where all missing values provided by an estimate of the posterior distribution. Second, we can continue to use one value of θ even after we have moved to a different value of θ .

4. MAP inference and sparse coding

4.3. inference refers to computing the probability distribution over one set of variables given another?

approach to learning with an approximate posterior.

3.3. E-step:

$$\text{Set } q_{h^{(i)}}(v) = p(h^{(i)} | v^{(i)}, \theta^{(t)})$$

$$\text{for all indices } i \text{ of the training examples } v^{(i)}$$

we want to train on (both batch and minibatch variants are valid), where $\theta^{(t)}$ is current parameter value.

missing values provided by an estimate of the posterior

to a different value of θ .

4. expectation maximization (EM).

4.3. inference refers to computing the probability distribution over

2.4. we can rearrange

convenient form:

$$L(v, \theta, q)$$

$$= E_{h^{(i)}} [log p(h^{(i)}, v)] + H(q)$$

examples $v^{(i)}$ we can think of inference as the procedure for finding the q that maximizes L . Exact

inference maximizes L perfectly by searching over a family of functions q that includes $p(h^{(i)}$) with respect to θ .

3.4. EM can be viewed as a coordinate ascent algorithm as maximizing L with respect to θ .

3.5. stochastic gradient ascent

3. expectation maximization (EM).

3.2. EM is not an inference, but rather an

Chapter 19 Approximate Inference

1. The challenge of inference usually refers to the difficult problem of computing $p(h|v)$ or taking expectations with respect to it.
2. Intractable inference problems in deep learning usually arise in interactions between latent variables in a structured graphical model.

2. Inference as optimization 2.2. Exact inference can be described as an optimization problem. Appropriate inference algorithms may be derived by approximating the underlying optimization problem.
- 2.3 evidence lower bound (ELBO), also called variational free energy.

3. expectation maximization $L(v, \theta, q) = log p(v; \theta) - D_{KL}(q(h|v) || p(h|v; \theta))$
- where q is an arbitrary probability distribution.

is the origin of the names

5.2 The core idea behind variational learning is that we can maximize L over a restricted family of distributions inferring the entire distribution.

5.3 The beauty of the variational approach is that we do not need to specify a specific parametric form for q . We specify how it should factorize, but then the optimization problem determines the optimal probability distribution within those factorization constraints.

5.4 Calculus of variations

MAP inference also us to find maximum a posteriori inference learn using a point estimate (MAP inference).

5.5 Mean field approach imposes the restriction that q is a factorial distribution:

$$q(h|v) = \prod_i q(h_i|v)$$

5.6 Structured variational inference imposes a graphical model structure we choose on q .

5.7 Learning algorithms based on

"variational learning" and "variational inference", though these names apply even when the latent variables are discrete and calculus of variations is not needed.

A typical way to do this is to introduce assumptions about how q factorizes.

5.8 Training sparse coding with maximum likelihood is intractable. Instead, we use MAP inference and learn parameters by maximizing the EIBI defined by the Dirac distribution around the MAP estimate of h .

5.9 MAP inference can be thought of as approximate inference.

5.10 MAP inference is commonly used in deep learning as both a feature extractor and a learning mechanism. It is primarily used for sparse coding. In some sense, we can think of MAP inference as approximate inference, because it does not provide the optimal q .

5.11 If $q(h|v) = \delta(h - \mu)$, // Dirac distribution then $\arg\max_q L(v, \theta, q)$ $\Leftrightarrow \mu^* = \arg\max_{\mu} p(h|v)$

5.12 MAP inference can be thought of as approximate inference.

- $J(y) = \int_X L(x, y(x), y'(x)) dx$ (1) functional derivative, also known as variational derivative.
- where x_1, x_2 are constants, $y(x)$ is twice continuously differentiable.
- $y'(x) = dy/dx$
- $L(x, y(x), y'(x))$ is twice continuously differentiable with respect to x, y, y' .
- If the functional $J(y)$ attains a local minimum at f , then
- $$\frac{\delta}{\delta f(x)} J = \frac{\partial L}{\partial f} - \frac{d}{dx} \frac{\partial L}{\partial f'} = 0$$
- (ii) normal distribution has the maximum entropy for fixed variance σ^2 .
- prove:
- $$= \frac{\partial}{\partial y} g(f(x), x)$$
- The Lagrangian functional:
- $$L(p) = \lambda_1 (Sp^2 dx - 1) + \lambda_2 (E[x] - \mu) + \lambda_3 (E[(x - \mu)^2] - \sigma^2) + H(p)$$
- 5.8' calculus of variations (2) we make a meanfield approximation:
- (1) we show how to apply variational inference to the binary sparse coding model.
- (ii) the derivative of the functional $J(f)$ with respect to the value of the function field approximation makes f at point x is denoted learning tractable.
- The ELBO L is tractable,
- $$\frac{\delta}{\delta f(x)} J$$
- and we use L as a replacement for the intractable maximum likelihood,
- (iii) for differentiable functions $g(y, x)$ with continuous derivatives, that
- $$\frac{\delta}{\delta f(x)} \int g(f(x), x) dx$$
- (iv) we use fixed point equations to estimate \hat{h} .
- (v) log-likelihood is intractable but $p(h|v)$ is a complicated distribution.
- (vi) log-likelihood We can resolve this difficulty by using variational inference and variational learning instead

implements an approximation $\hat{f}(v, \theta)$ offsets the learning process, and this in turn affects the accuracy of the inference algorithm.

6.3 one possible explanation for the role of dream sleep is that dream could provide the negative phase samples that Monte Carlo training algorithms use to approximate the negative gradient of the log partition function of undirected models. Another possible explanation is that it is providing samples from $p(h, v)$ which can be used to train an inference network to predict h given v . Dreaming may also serve reinforcement learning rather than probabilistic modeling, by sampling synthetic experiences from the animal's transition model, on which to train the animal's policy.

5.9 continuous latent variables

If we make the mean field approximation

$$q(h|v) = \prod_i q(h_i|v)$$

(ii) The training algorithm tends to adapt the model in a way that makes the approximating assumptions underlying the approximate inference algorithm $q(h|v)$ may be obtained become more true.

6. learned approximate inference 6.2 Inference can be thought of as an optimization procedure that increases the value of L . Once we think of the multi-step iterative optimization process as just being a function that maps an input v to an approximate distribution

$q^* = \arg\max_L(v, q)$, we can approximate it with (i) using approximate as a neural network that

To minimize the Lagrangian with respect to p , we set the functional derivatives equal to 0:

$$\begin{aligned} \nabla_x \cdot \frac{\delta}{\delta p} L \\ = \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - \\ - \log(p(x)) \\ = 0 \end{aligned}$$

$\Rightarrow p(x) = \exp(\lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2)$

We are free to choose any λ values, because the gradient of the Lagrangian with respect to the λ variables is zero so long as the constraints are satisfied.

We may set $\lambda_1 = -\log(p(x))$, $\lambda_2 = 0$, and $\lambda_3 = -\frac{1}{2\sigma^2}$

5.10 interactions between learning and inference

part of a learning algorithm

training of deep architectures. 2.5 $p(h|v) = \frac{p(h,v)}{p(v)}$

The introduction of DBNs in 2006 began the current deep learning renaissance.

3.2 DBNs are generative models with several layers of latent variables. The latent variables 2.6 $p(h_j=1|v) = \sigma(c_j + v^T W_{j,j})$ are typically binary, while the visible units may be binary. Similarly, $p(v|h) = \prod_{i=1}^{n_v} \sigma(c_i + h^T W_{i,h})$, or real. There are no intra-layer connections. The connections between the top two layers are undirected. The connections between all other layers are directed, with the arrows pointed toward the layer that is closest to the data.

3.3 A DBN with L hidden layers:

$$p(h^{(1)}, h^{(L-1)}) \propto \exp(c^{(1)}h^{(1)} + \dots + c^{(L-1)}h^{(L-1)} + h^{(L-1)}W^{(L-1)}h^{(L)})$$

makes Boltzmann machine learning somewhat biologically plausible.

Chapter 20 deep generative models

1. Boltzmann machines

1.2 Boltzmann machines were originally introduced as a general "connectionistic" approach to learning arbitrary probability distribution over binary vectors.

2. RBMs

2.2 RBMs are undirected probabilistic graphical models containing a layer of observable variables and a single layer of latent variables.

i. $p(h|v)$ is factorial.

ii. $p(h_j=1|v) = \sigma(c_j + v^T W_{j,j})$

iii. $p(h|v) = \prod_{j=1}^{n_h} \sigma(c_j + h^T W_{j,h})$

(See RBM in Chapter 16.)

when all the variables are observed

$E(x) = -x^T U x - b^T x,$

$E(v, h) = -v^T R v - v^T N h - h^T S h$

where we decompose the units x into two subsets: the visible units v and the latent units h .

1.3 $p(x) = \frac{\exp(-E(x))}{Z}$

1.4 a Boltzmann machine with hidden units becomes a universal approximator of probability mass functions over discrete variable

2.4. $p(h|v)$ and $p(v|h)$ are factorial and relatively simple functions over discrete variable

3.1 DBNs were one of the first non-convolutional models to compute and to sample

3.2 deep belief networks (DBNs)

3.3 The learning rule of Boltzmann machine is "local," which

4.7 The use of proper mean field allows the approximate inference procedure for DBM to capture the influence of top-down feedback interactions. This makes DBMs interesting from the point of view of neuroscience.

4.8 Sampling from DBMs is relatively difficult.

4.9 Learning in the DBM must confront both the challenge of an intractable partition function and the challenge of an intractable posterior distribution.

4.10 DBMs are typically trained using SMN (variational SMN). Many of the other techniques described in chapter 18 are not applicable.

on one side and even layers on the other.

4.5 The bipartite structure makes Gibbs sampling in a DBM efficient. The naive approach to Gibbs sampling is to update only one variable at a time. Gibbs sampling in a DBM can be divided into two blocks of updates, one including all even layers (including the visible layer) and the other including all odd layers.

4.8 Sampling from DBMs is relatively difficult.

4.9 Learning in the DBM must confront both the challenge of an intractable partition function and the challenge of an intractable posterior distribution.

4.10 DBMs are typically trained using SMN (variational SMN). Many of the other techniques described in chapter 18 are not applicable.

- 3.6 We use a greedy layer-wise procedure to train a DBM.
4. deep Boltzmann machines (DBMs).
- 4.2. DBM is another kind of deep generative model. Unlike DBV, it is an entirely undirected model. $\forall i$.
- Unlike RBM, DBM has several layers of latent variables. But like RBM, within each layer, each of the variables are mutually independent, conditioned on the variables in the neighboring layers.
- 3.4 A DBN with only one hidden layer is just an RBM.
- 3.5 To generate a sample from a DBN, we first run several steps of Gibbs sampling on the top two hidden layers, to omit the bias parameters, $E(v, h^{(1)}, h^{(2)}, \theta)$
- 4.3 In the case of a DBM with three hidden layers, the top two hidden layers, is simpler for DBM.
- The simplicity of this posterior distribution allows richer approximations
- 4.4 the DBM layer can be organized into a bipartite graph, with odd layers
- 3.6 We use a greedy layer-wise procedure to train a DBM.
- $p(h_i^{(k)} = 1 | h^{(k+1)}) = \sigma(b_i^{(k)} + W_{:,i}^{(k+1)\top} h^{(k+1)})$, $\forall i, \forall k \in [1, m], l=2$.
- $p(v_i=1 | h^{(0)}) = \sigma(b_i^{(0)} + W_{:,i}^{(0)\top} h^{(0)})$,
- In the case of real-valued visible units, substitute $v \sim N(v; b^{(0)} + W^{(0)\top} h^{(0)}, \beta^{-1})$ with β diagonal for tractability.
- 3.4 A DBN with only one hidden layer is just an RBM.
- 3.5 To generate a sample from a DBN, we first run several steps of Gibbs sampling on the top two hidden layers, to omit the bias parameters, the top two hidden layers, then use a single pass of ancestral sampling through the rest of the model to draw a sample from the visible units.

the canonical energy model back-propagation computes the exact gradient of the loss. The disadvantage is that it does not provide a way to optimize the log-likelihood, but rather a heuristic approximation of the generalized pseudolikelihood.

5.4 Since the Gaussian RBM

only models the conditional mean of the input given the hidden units, it cannot capture conditional covariance information. Alternative

models have been proposed that attempt to better account connections to dropout, for the covariance of real-valued data. These models include for real-valued data the mean and covariance

RBMs (mRBMs), the mean product of t-distribution (mPT) model and the spike and slab RBMs (SSRBMs)

$\Rightarrow p(v|h) = N(v; Wh, \beta^{-1})$, where h is binary,

v is real-valued.

5.5 convolutional Boltzmann machines.

with MC estimates of the gradient.

The disadvantage is that MP-RBM works by viewing pretraining of a DBN.

It does not provide a way to optimize the log-likelihood, as defining a family of recurrent networks for approximately solving every possible inference problem.

Rather than training the model to maximize the likelihood, the model is trained to make each recurrent network obtain an accurate answer to the corresponding inference problem.

$$-B^T h$$

$E(v|h) = \frac{1}{2} v^T (B \otimes v) - (v \otimes B^T) h$

v is real-valued.

5.6 convolutional Boltzmann machines.

captures the exact gradient of the loss. The disadvantage is that MP-RBM works by viewing pretraining of a DBN.

There are two main ways to resolve the joint training problem of DBM. The first is centered DBM. The second way to jointly train a DBM is to use a multi-prediction DBM (MP-DBM).

4.13 centered DBM reparametrizes the model in order to make the Hessian of the cost function better-conditioned at the beginning of the learning process. This yields a model that can be trained without a greedy layer-wise pretraining stage.

4.14 MP-DBM uses an alternative training criterion that allows the use of the back-propagation algorithm

as it is really in order to avoid the problems that

captures the exact gradient of the loss. The disadvantage is that MP-DBM works by viewing pretraining of a DBN.

4.12 There are two main ways to resolve the joint training problem of DBM. The first is centered DBM. The second way to jointly train a DBM is to use a multi-prediction DBM (MP-DBM).

4.13 centered DBM reparametrizes the model in order to make the Hessian of the cost function better-conditioned at the beginning of the learning process. This yields a model that can be trained without a greedy layer-wise pretraining stage.

4.14 MP-DBM uses an alternative training criterion that allows the use of the back-propagation algorithm

as it is really in order to avoid the problems that

derivatives of y with respect to the parameters of its distribution, μ and σ^2 .

9.3 We can express any probability distribution of the form $p(y; \theta)$ or $p(y|x; \theta)$ as $p(y|\theta)$, where θ is a variable containing both parameters θ , and if applicable, the inputs x .

8.4 The Boltzmann machine 9.4 Given a value y sampled from distribution $p(y|w)$, where w may in turn be a function of other variables, we can rewrite $y \sim p(y|w)$ as $y = f(z; w)$, where z is a source of randomness. We may then compute the derivatives of y with respect to w applied to random operations using back-propagation algorithm, so long as f is continuous and differentiable almost everywhere. Crucially, w must not be a function of z , and z must not be a function of w . This technique,

as powerful supervised learners as MLPs, at least using existing methodology.

7.4 Another sequence modeling task is to model the distribution over sequences of musical notes feature map that can scale in size proportional to the size of the input image.

7.5 Boltzmann machines for sequence modeling can be applied to this task. 7.6 Boltzmann machines for structured or sequential output

8. other Boltzmann machines 8.2 Boltzmann machines may be extended with different training criteria. We have focused on Boltzmann machines trained to approximately maximize the generative criterion $\log p(y)$. It is also possible to train discriminative RBMs that aim to maximize $\log p(y|v)$ instead. Unfortunately, RBMs do not seem to be able to take the

RBM modeling problem [10]. 7.3 Boltzmann machines that for small m can be applied to use probabilistic max pooling may accept variable-sized this task.

7.4 Another sequence modeling task is to model the distribution over sequences of musical notes feature map that can scale in size proportional to the size of the input image.

7.5 Boltzmann machines for sequence modeling can be applied to this task. 7.6 Boltzmann machines for structured or sequential output

8. other Boltzmann machines 8.2 Boltzmann machines may be extended with different training criteria. We have focused on Boltzmann machines trained to approximately maximize the generative criterion $\log p(y)$. It is also possible to train discriminative RBMs that aim to maximize $\log p(y|v)$ instead. Unfortunately, RBMs do not seem to be able to take the

7.4 Another sequence modeling task is to model the distribution over sequences of musical notes feature map that can scale in size proportional to the size of the input image.

7.5 Boltzmann machines for sequence modeling can be applied to this task. 7.6 Boltzmann machines for structured or sequential output

8. other Boltzmann machines 8.2 Boltzmann machines may be extended with different training criteria. We have focused on Boltzmann machines trained to approximately maximize the generative criterion $\log p(y)$. It is also possible to train discriminative RBMs that aim to maximize $\log p(y|v)$ instead. Unfortunately, RBMs do not seem to be able to take the

samples x or to distributions 10.3 DBNs are partially directed model.

over samples x using a differentiable function $g(z; \theta^g)$

which is typically represented by a neural network. This model class includes VAEs,

10.5 sigmoid belief networks QANs, and techniques that are a simple form of directed graphical model with a specific kind of conditional probability

in isolation.

10.7 variational autoencoder (VAE) distribution. We can think of (i) VAE is a directed model that a sigmoid belief network as uses learned approximate if one having a vector of binary and can be trained purely with states s , with each element of the state influenced by its ancestors:

the model, the VAE first draws $p(s_i) = \sigma(\sum_j W_{ji} s_j + b_i)$

a sample \bar{s} from the code distribution $P_{model}(z)$. The sample

is then run through a differentiable generator, transforms samples generator network $g(z)$. Finally, of latent variables \bar{x} to x is sampled from a distribution

9.7 $E_z[J(y)] = \bar{y} J(\bar{y}) p(y)$ is often called the reparametrization trick, stochastic back-propagation or perturbation analysis.

// assume J does not reference y . If we wish to back-propagate w directly, here we can relax the assumption.

$= \bar{y} J(\bar{y}) \frac{\partial p(y)}{\partial w}$ that produces discrete-valued samples, it may still be possible to estimate a gradient on w , using reinforcement learning algorithms such as // an unbiased Monte Carlo estimator variants of the REINFORCE algorithm.

9.6 The core idea of the REINFORCE algorithm is that even though $J(f(z; \omega))$ is a step function with useless derivatives, the expected learning community, within the smaller deep learning community, they have until roughly 2013 been overshadowed by undirected gradient descent models such as the RBM.

(ii) The ~~tech~~ networks are trained the training data and samples encoder and decoder with a technique called moment drawn from the generator: (iv) The VAE approach is matching.

The discriminator emits a probability among the state of the art (iii) generative moment matching value given by $d(x; \theta^d)$, approaches to generative indicating the probability that modeling.

it is a real training example (v) VAE is an excellent rather than a fake sample manifold learning algorithm. called maximum mean discrepancy drawn from the model.

or MMD. The cost function 10.8 generative adversarial networks (GANs) measures the error in the first moments in an infinite-dimensional networks space, using an implicit mapping (i) generative moment matching modeling approach based on that produces the parameters networks are another form of generative model based on a kernel function in order to make computations on infinite-dimensional vectors tractable. The MMD cost is zero if and only if the two distributions being compared are equal.

10.9 generative moment matching (GANs)

(ii) GANs are another generative (also called an inference network) theoretic scenario in which the generator network must samples of z drawn from compete against an adversary. $q(z|x) = q(z; f(x; \theta))$ The generator network directly in order to obtain a grade produces samples $x=g(z; \theta)$. with respect to θ . Learning with VAEs nor a discriminator Its adversary, the discriminator then consists solely of networks as used with GANs. attempts to distinguish maximizing L with respect between samples drawn from to the parameters of the

$P_{\text{model}}(x; g(z)) = P_{\text{model}}(x|z)$. However, during training, the approximate inference network (encoder) $q(z|x)$ is used to obtain z and $P_{\text{model}}(x|z)$ is then viewed as a decoder network.

deriving autoencoders:

- (i) starting from the previous state x , the new parametrization can inject corruption noise, sampling \tilde{x} from $C(\tilde{x}|x)$.

(ii) encode \tilde{x} into $h = f(\tilde{x})$.
(iii) decode h to obtain the parameters principle.

$$\omega = g(h) \text{ or } p(x|w=g(h)) = p(x|\omega)$$

- (iv) sample the next state x from density estimator (NVAE)

$$p(x|w=g(h)) = p(x|\tilde{x})$$

12. generative stochastic networks
(A SNS)

12.2 ASNs are generalizations of denoising autoencoders that include latent variables h in the generative Markov chain, in addition to the visible variables x .

13. generative models hold the promise to provide AI systems with a framework for all of the many different intuitive concepts they need to understand, and the ability to reason about these concepts in the face of uncertainty.

- increased as much as needed. Each $p(x|x_{t-1}, \dots, x_1)$ is parameterized as a linear directed probabilistic model. (ii) linear auto-regressive nets are with no latent random variable. The conditional probability distributions in these models are represented by neural networks. They decompose a joint probability over the observed variables using the chain rule of probability to obtain a product of conditionals of the form $p(x_t|x_{t-1}, \dots, x_1)$. Such models have been called fully-visible graphical model as logistic auto-regressive networks but Bayes networks (PVBMs).
- 10.10 auto-regressive networks
(i) linear auto-regressive nets also improve generalization by introducing a parameter sharing and feature sharing
- (ii) encode \tilde{x} into $h = f(\tilde{x})$.
generalization of linear classification methods to generative modeling.
- (iii) (i) NVAE (ii) linear classifiers
(i) NVAE is a very successful the model itself does not offer a way of increasing its capacity.
- The parameters of the hidden units of different group are shared.
- (ii) forward propagation in a NVAE model loosely resembles the computations performed in mean field inference to fill in missing inputs in an RNN within that graphical model is the simplest form of the new parametrization auto-regressive network which has more powerful in the sense no hidden units and no sharing that its capacity can be of parameters or features.

4.2 vector space $(V, +, \cdot)$
over a field F .

$\Leftrightarrow (V, +, \cdot)$ is left R -module,
where F is the ring.

5. algebra A over a field K

$\Leftrightarrow A$ is a vector space over K ,
with a binary operation $\cdot : A \times A \rightarrow A$,
satisfying:

$$\begin{aligned} \forall x, y, z \in A, \quad a, b \in K, \\ \text{① } (x+y) \cdot z = x \cdot z + y \cdot z \\ \text{② } x \cdot (y+z) = x \cdot y + x \cdot z \\ \text{③ } (a \cdot x) \cdot y = a \cdot (x \cdot y) \end{aligned}$$

6. algebraic structure is a set

with one or more operations
defined on it that satisfies
a list of axioms.

6) examples: groups, rings, fields,
lattices, vector spaces,
modules, algebras.

3.2 lattice (L, \vee, \wedge)

\Leftrightarrow ① (L, \vee) is commutative semigroup.
 \Leftrightarrow ring without assuming

② (L, \wedge) is commutative semigroup.
the existence of multiplicative identity.

③ absorption law.

$$a \vee (a \wedge b) = a$$

$$a \wedge (a \vee b) = a$$

$$\forall a, b \in L.$$

3.3 bounded lattice

$$(L, \vee, \wedge, 0, 1)$$

④ distributivity of multiplication over addition.

$\Leftrightarrow (L, \vee, \wedge)$ is lattice
0 is identity element for \vee .

1 is identity element for \wedge .

$$\forall a, b, c \in F.$$

4. module:

R is a ring, left R -module $(M, +, \cdot)$

$\Leftrightarrow (M, +)$ is an abelian group.

$$+: R \times M \rightarrow M$$

satisfying, $\forall r, s \in R, x, y \in M$

$$\begin{aligned} \text{① } r \cdot (x+y) &= r \cdot x + r \cdot y \\ \text{② } (r+s) \cdot x &= r \cdot x + s \cdot x \\ \text{③ } (r \cdot s) \cdot x &= r \cdot (s \cdot x) \\ \text{④ } 1_R \cdot x &= x, \quad 1_R \text{ is multiplicative identity of } R. \end{aligned}$$

2.2 Algebraic structures

set S

1. magma \Leftrightarrow set with a binary operation
 $\cdot : S \times S \rightarrow S$.

\Leftrightarrow semigroup \Leftrightarrow associative magma

\Leftrightarrow monoid \Leftrightarrow semigroup with identity

\Leftrightarrow group \Leftrightarrow monoid with existence of inverses

\Leftrightarrow abelian group \Leftrightarrow group
except 0 doesn't have multiplicative inverse.

\Leftrightarrow

2. ring $(R, +, \cdot)$

\Leftrightarrow $(R, +)$ is abelian group

⑤ (R, \cdot) is monoid.

③ multiplication is distributive with respect to addition

$$\begin{aligned} a \cdot (b+c) &= (a \cdot b) + (a \cdot c) \\ (b+c) \cdot a &= (b \cdot a) + (c \cdot a) \end{aligned}$$

$\forall a, b, c \in R.$

7. $f: \mathbb{R}^m \rightarrow \mathbb{R}$

$f'(x)=0, f''(x)$ positive definite \Rightarrow
 x is local minimum.

proof: $\forall u$ is a unit direction,

$$\frac{\partial^2 f}{\partial u^2} = u^T (f''(x))^T u > 0$$

($\because f''$ pd $\Rightarrow (f'')^T$ pd)

$$\therefore \frac{\partial^2 f}{\partial u^2} = f'(x)^T u = 0.$$

$\therefore x$ is local minimum along u .
 (get from 6).

$$\because u$$
 is arbitrary direction
 $\therefore f(x)$ can be approximated using $h(x)$.

$\therefore x$ is local minimum of f .

7.2. $f'(x)=0, f''(x)$ negative definite

$\Rightarrow x$ is local maximum.

7.3. $f'(x)=0, f''(x)$ symmetric,
 $\lambda_1 > 0, \lambda_2 < 0$ are eigenvalues of $f''(x)$,
 $\Rightarrow x$ is saddle point.

5. f'' measures curvature.

$$\text{if } f(x) = f(a) + f'(a)(x-a) + f''(x)(x-a)^2$$

$$+ \frac{f'''(a)}{2}(x-a)^3 + o((x-a)^3)$$

$\therefore f''(a) > 0$

$$\frac{\partial^2 f}{\partial x^2} = u^T (f''(x))^T u > 0$$

set $h(x) = f(x) - f(a) - f'(a)(x-a) + \frac{f''(a)}{2}(x-a)^2$

$\therefore h(x), h'(x)$ is \wedge ,

$\therefore f''(a) < 0, h(x)$ is \wedge ,

$\therefore f''(a) = 0, h(x)$ is \sim

$\therefore f(x) = h(x) + o((x-a)^2)$

$\therefore f(x)$ can be approximated using $h(x)$.

$\therefore f''(a)$ also measures curvature of f at a .

7.2. $f'(x)=0, f''(x)$ negative definite

$\therefore f(x)$ is local minimum.

7.3. $f'(x)=0, f''(x)$ symmetric,

$\lambda_1 > 0, \lambda_2 < 0$ are eigenvalues of $f''(x)$,

$\Rightarrow x$ is saddle point.

$\therefore x$ is local maximum.
 (we can see that from 5).

Chapter 4.

1. $f(x): \mathbb{R}^2 \rightarrow \mathbb{R}$

$\|u\|=1$, u is a direction, and fixed.

H is called Hessian matrix of f . $\Rightarrow \frac{\partial^2 f}{\partial u^2} = \lim_{\alpha \rightarrow 0} \frac{f(x+\alpha u) - f(x)}{\alpha}$

$$= u^T \nabla^2 f(x) = (f'(x))^T u$$

where $f'(x) := \nabla f(x)$, $\in \mathbb{R}^2$

2. $f(x): \mathbb{R}^m \rightarrow \mathbb{R}^2$, $\|u\|=1, u$ is direction

$\therefore f(x) = (f_1(x), f_2(x))$

$$\frac{\partial^2 f}{\partial u^2} := \left(\frac{\partial^2 f_i}{\partial x_j^2} \right)_{i,j} = \left(\frac{\partial f_i}{\partial x_j} \right)_U = J_u$$

$\therefore f(x) = (f_1(x), f_2(x))$, is Jacobian of f .

$f: R \rightarrow R$, is k times differentiable, where $J := \begin{pmatrix} (Df_1)^T \\ (Df_2)^T \end{pmatrix}$, is Jacobian of f .

at a , $\Rightarrow \exists h_a: R \rightarrow R$.

such that

$f(x) = f(a) + f'(a)(x-a) + \dots$ then $\frac{\partial^2 f}{\partial u^2} = J^T \in \mathbb{R}^{m \times 2}$

3. $f(x): \mathbb{R}^2 \rightarrow \mathbb{R}$

$\therefore x$ is local minimum.

$\therefore f(x) = f'(x) = 0$,
 $\text{and } \lim_{x \rightarrow a} h_k(x) = 0$

$\therefore g(x): \mathbb{R}^2 \rightarrow \mathbb{R}^2$

set H is Jacobian of g .

proof:

$$f(x) = f(x_0) + (x-x_0)f'(x_0) + \frac{1}{2}(x-x_0)^2 f''(x_0) + o((x-x_0)^2)$$

$$8. f(x) = kx + o(x), k > 0$$

$$\Rightarrow \exists \delta > 0, \forall x \in (0, \delta), \\ \text{we have } f(x) > 0.$$

$$\text{Let } x = x_0 + \alpha, \because f'(x_0) = 0$$

$$\therefore f(x_0 + \alpha) = f(x_0) + \frac{1}{2}\alpha^2 f''(x_0) + o(\alpha^2)$$

$$1. f(x_0 + \alpha) - g(x_0 + \alpha)$$

$$= \frac{1}{2}\alpha^2(f''(x_0) - g''(x_0)) + o(\alpha^2)$$

proof: from 9, we know,

$$\exists \delta > 0, \forall x \in (0, \delta), \\ |o(x)| < \frac{k}{2}|x|$$

$$\therefore o(x) > -\frac{k}{2}x, \forall x \in (0, \delta)$$

$$\therefore f(x) = kx + o(x) > \frac{k}{2}x > 0, \\ \forall x \in (0, \delta).$$

$$9. f(x) = o(g(x)), x \rightarrow a$$

$$\Leftrightarrow \lim_{x \rightarrow a} \frac{f(x)}{g(x)} = 0$$

$$\Leftrightarrow \forall \varepsilon > 0, \exists \delta > 0, \forall x \in (a, \delta),$$

$$\left| \frac{f(x)}{g(x)} \right| < \varepsilon$$

$$\Leftrightarrow \forall \varepsilon > 0, \exists \delta > 0, \forall x \in (a, \delta),$$

$$\text{we have } |f(x)| < \varepsilon |g(x)|, g(x) \neq 0.$$

$$11. f: \mathbb{R}^m \rightarrow \mathbb{R}, f', f'' \text{ exist},$$

x_0 is local minimum, $\lambda_1 > \lambda_2$

$$\Rightarrow f(x_0 + \alpha v_1) > f(x_0 + \alpha v_2) \\ \exists \delta > 0, \forall \alpha \in O(0, \delta),$$

proof: $\because x_0$ is local minimum.

f'' is symmetric,

λ_1, λ_2 is eigenvalues of f'' ,

v_1, v_2 is eigenvectors of f'' .

$$\because \lambda_1 > \lambda_2, \therefore \frac{\partial^2 f}{\partial v_1^2} > \frac{\partial^2 f}{\partial v_2^2}$$

from 10, we get the conclusion.

$$10. f: \mathbb{R} \rightarrow \mathbb{R}, g: \mathbb{R} \rightarrow \mathbb{R},$$

at x_0, f', f'', g', g'' exist.

x_0 is local minimum of f, g .

$$f''(x_0) > g''(x_0), f(x_0) = g(x_0)$$

$$\Rightarrow \exists \delta > 0, \forall \alpha \in O(0, \delta), f(x_0 + \alpha) > g(x_0 + \alpha)$$

5. PCA

prof: $E[\hat{\sigma}_m^2]$

5.2 PCA finds a representation $\tilde{x} = x^T W$, where $\text{Var}(\tilde{x})$ is diagonal.

Here W comes from SVD, $x = U \Sigma V^T$.

$$\text{prof: } \text{Var}(\tilde{x}) = \frac{1}{m} \tilde{x}^T \tilde{x}$$

3. $\{x^{(1)}, \dots, x^{(m)}\}$ are points,

$$x^{(i)} \in \mathbb{R}^n, \forall i, \text{ let } X \in \mathbb{R}^{m \times n}, X_{i,:} = x^{(i)T}$$

$$= \frac{1}{m} W^T X^T X W \quad \text{if } E(x^{(i)}) = 0, \forall i,$$

$$= \frac{1}{m} \Sigma^2 \quad \text{then unbiased sample covariance matrix is } \frac{1}{m} X^T X.$$

5.3 W in PCA is a rotation of x that aligns the principal axes of variance with the basis of new representation space associated with 2.

4.2 three criteria of a simple representation.

1. $\hat{\sigma}_m^2$ is biased.

5.4. PCA is an unsupervised learning algorithm that learns a representation of data. The representation has a lower dimension. The elements of the representation are no linear correlation with each other.

$$\therefore EA = \frac{m}{2} E(x^{(i)})^2$$

$$= \frac{m}{2} (\text{Var}(x^{(i)}) + E(x^{(i)}))^2$$

$$= m\sigma^2 + m\mu^2$$

$$EB = mE(\hat{\mu}^2)$$

$$= mE\left(\left(\frac{\sum x^{(i)}}{m}\right)^2\right)$$

1. proved.

2. Estimators of variance of Gaussian distribution.

$$= \frac{1}{m} (\text{Var}(\Sigma x^{(i)}) + (E(\Sigma x^{(i)}))^2)$$

$$= \frac{1}{m} (\Sigma \text{Var}(x^{(i)}) + (\Sigma E(x^{(i)}))^2)$$

$$= \frac{1}{m} (m\sigma^2 + m^2\mu^2) \quad \text{2.2 sample variance of } \sigma^2$$

$$\hat{\sigma}_m^2 = \frac{1}{m} \sum (x^{(i)} - \hat{\mu}_m)^2$$

where $\hat{\mu}_m = \frac{1}{m} \sum x^{(i)}$ is sample mean $\{x^{(1)}, \dots, x^{(m)}\}$ are i.i.d. according to $p(x^{(i)}) = N(x^{(i)}; \mu, \sigma^2)$ $\forall i \in \{1, \dots, m\}$.

$\Rightarrow \hat{\sigma}_m^2$ is biased estimator.

1. $\hat{\sigma}_m^2$ is biased.

2.3 unbiased sample variance estimator.

$$\hat{\sigma}_m^2 = \frac{1}{m-1} \sum (x^{(i)} - \hat{\mu}_m)^2$$

$\Rightarrow \hat{\sigma}_m^2$ is unbiased.

$$= \frac{1}{m} E(A - B)$$

Chapte5.

1. A is real symmetric matrix

$$f(x) = x^T A x, \|x\| = 1.$$

$\Rightarrow \min_x f(x)$ is minimum eigenvalue of A .

$\max_x f(x)$ is maximum eigenvalue of A .

2. Estimators of variance of Gaussian

distribution.

$$= \frac{1}{m} (\text{Var}(\Sigma x^{(i)}) + (E(\Sigma x^{(i)}))^2)$$

$$= \frac{1}{m} (\Sigma \text{Var}(x^{(i)}) + (\Sigma E(x^{(i)}))^2)$$

$$= \frac{1}{m} (m\sigma^2 + m^2\mu^2) \quad \text{2.2 sample variance of } \sigma^2$$

$$\hat{\sigma}_m^2 = \frac{1}{m} \sum (x^{(i)} - \hat{\mu}_m)^2$$

where $\hat{\mu}_m = \frac{1}{m} \sum x^{(i)}$ is sample mean $\{x^{(1)}, \dots, x^{(m)}\}$ are i.i.d. according to $p(x^{(i)}) = N(x^{(i)}; \mu, \sigma^2)$ $\forall i \in \{1, \dots, m\}$.

$\Rightarrow \hat{\sigma}_m^2$ is biased estimator.

prof: $E[\hat{\sigma}_m^2] = \frac{1}{m} \left[\frac{m}{2} ((x^{(1)})^2 - \hat{\mu}_m)^2 \right]$

$$= \frac{1}{m} E\left[\frac{m}{2} (x^{(1)})^2 + \hat{\mu}_m^2 - 2\hat{\mu}_m x^{(1)}\right]$$

$$= \frac{1}{m} E\left[\frac{m}{2} (x^{(1)})^2 - m\hat{\mu}^2\right]$$

$$= \frac{1}{m} E(A - B)$$