

1 - MODERN HASHING

CONSISTENT HASHING AND CMS

Consistent Hashing Very useful when we have a hashstable distributed across multiple machines (e.g. shared cache). Hash all objects, but also hash all names of servers (to the same range), store objects in first server hashed to their right around the circle. Can reduce variance with virtual replication (with copies proportional to storage capacity). When new server is added, place it on the circle and rehash the things between it and the preceding bucket. In $\log(n)$ time, use balanced binary search tree to look up in $\log(n)$ time, where n is the number of servers.

CMS There is no algorithm that exactly solves Heavy Hitters in one pass while using sublinear amount of auxiliary space. Approximate: require every value that occurs at least $\frac{n}{k}$ times repeated, but also that every value in the list occurs at least $\frac{n}{k} - \epsilon n$ times in A (some false positives). CMS data structure has two operations, Inc(x) and Count(x). $I \times b$ table of I hash functions each with b buckets. For Count(x), return min value observed, for Inc(x) increment list of set of minimum values. Bloom filter is CMS but with only binary entries - entry checking is "are ALL the buckets 1". Space required is bl counters where $b = \frac{n}{k}$ and $l = \ln \frac{1}{\delta}$, where δ is error probability. If n the length of the stream is not known in advance, keep track of running size of stream m and store elements in a min-heap that occur more than $\frac{m}{k}$ times. Prune minimum element if it falls below $\frac{m}{k}$. Heap contains at most $2k$ elements, so maintaining heap requires an extra $O(\log k)$ work per array entry.

2 - DATA WITH DISTANCES

Jaccard $J(S, T) = \frac{|S \cap T|}{|S \cup T|}$

$$L_p ||x - y||_p = (\sum_{i=1}^n |x_i - y_i|^p)^{1/p}$$

Voronoi diagram Given points X in \mathbb{R}^d , and some metric D , the Voronoi diagram partitions space into regions (the set of all points that are closest to a single datapoint). Concretely, for $x \in X$,

$part(x) = \{y \in \mathbb{R}^d \text{ s.t. } D(x, y) = \min_{x' \in X} D(x', y)\}$
 k -tree Given a set $S = \{x_1, x_2, \dots, x_n\}$. Pick a dimension / coordinate i . Compute a median $\{x_i\}$. Partition: $S_1 = \{x \in S | x_i < m\}$, and $S_2 = \{x \in S | x_i \geq m\}$. Recurse on S_1 and S_2 and store which dimension we are looking at. Runtime is logarithmic in number of points, and exponential in dimension of points.

THE KISSING NUMBER. How many identical spheres can you place around a sphere such that all of them are touching the center sphere?

MinHash Consider Jaccard similarity, defined earlier. We develop the "MinHash" technique that we can use to reduce dimensionality. Pick a uniformly random ordering of the universe U . Map set S to $f(S) =$ "min" element of S . For any sets S and T , we have $Pr(f(S) = f(T)) = J(S, T)$.

MinHash Minhash preserves Jaccard similarities in expectation, like JL preserves L_2 in expectation. MinHash is an unbiased estimator of Jaccard Similarity. If we want an accurate estimate up to $\pm \epsilon$ of all n choose 2 Jaccard similarities, use $O(\epsilon^{-2} \log n)$ independent estimators. Locality sensitive hashing is an extension that maps close things to close buckets (so you can search nearby buckets to remove near-duplicates, which is not possible with a hash). **KD Trees** Perform well in < 20 dimensions, or if number of points is at least 2^d . Build a binary search tree that partitions space; edges of tree correspond to subsets of space, and each node v has the index of some dimension i_v and a value m_v . Generally choose random dimension first of dimen choose m_v to be median of that dimension's coordinates. Given a node v , first go down the tree and find the leaf in which v would end up. Then go back up the tree, at each junction, asking "Is it possible that the closest point to v would have ended up down the other path?". In low dimensions, answer is almost always no, but in high dimensions we end up exploring a lot (exponential in number of dimensions due to curse of dimensionality - hence bad performance).

JL transformation First, note that normals are closed. Let $X \sim N(\mu_1, \sigma_1^2)$, $Y \sim N(\mu_2, \sigma_2^2)$. Then $X + Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$. Note that this is fairly unique to Gaussians - this isn't true for most distributions. For any two vectors $x, y \in \mathbb{R}^d$, we have $E[|(f(x) - f(y))^2|] = E[||x - y||^2]$ Proof is easy:

$$E[|(f(x) - f(y))^2|] = E\left\{\left|\sum_{i=1}^d r_i x_i - \sum_{i=1}^d r_i y_i\right|^2\right\}$$
$$= E\left\{\left|\sum_{i=1}^d r_i (x_i - y_i)\right|^2\right\}$$
$$= E\left\{\left\langle 0, \sum_{i=1}^d (x_i - y_i)^2 \right\rangle\right\}$$
$$= \text{Var}\left\{\left\langle 0, \sum_{i=1}^d (x_i - y_i)^2 \right\rangle\right\}$$
$$= \sum_{i=1}^d (x_i - y_i)^2.$$

Note: If you repeat $t = \frac{1}{\epsilon^2 \log n}$ times, then with high probability, all $\binom{n}{2}$ pairwise distances are preserved to within a factor of $1 \pm \epsilon$. This transformation is referred to as the Johnson-Lindenstrauss transformation. **JL Transform** Gaussians are closed under addition! $X_1 + X_2$ has distribution $N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$. Variances add for independent random variables, expectations add for any variables. For a set of n points in k dimensions, to preserve all n choose 2 interpoint Euclidean distances up to $\pm \epsilon$ factor, set $d = O(\epsilon^{-2} \log n)$. JL transform sends $d \rightarrow \frac{1}{\epsilon^2} A x$ - don't forget scaling factor of d^{-1} ! Preserves L_2 distances in reduced space. A is a $d \times k$ matrix. Can also just use ± 1 entries.

MinHash Minhash preserves Jaccard similarities in expectation, like JL preserves L_2 in expectation. MinHash is an unbiased estimator of Jaccard Similarity. If we want an accurate estimate up to $\pm \epsilon$ of all n choose 2 Jaccard similarities, use $O(\epsilon^{-2} \log n)$ independent estimates.

Locality sensitive hashing is an extension that maps close things to close buckets (so you can search nearby buckets to remove near-duplicates, which is not possible with a hash).

3 - GENERALIZATION

Consider the binary classification setting. We can broadly define it as follows: Suppose we have some datapoints $x_1, \dots, x_n \in \mathbb{R}^d$. Known distribution D on \mathbb{R}^d . Ground truth label function $f: \mathbb{R}^d \rightarrow \{0, 1\}$. **PROBLEM.** Given x_1, \dots, x_n drawn independently from D , and labels $f(x_1), \dots, f(x_n)$, our goal is to output $\hat{f}: \mathbb{R}^d \rightarrow \{0, 1\}$ such that $\hat{f}_n \approx f^n$. Namely, we want the generalization error to be low, defined this way: generalization error(g) = $Pr_{x \sim D} [g(x) \neq f(x)]$. Can also define the training error as the fraction of the training points on which g disagrees with the true labelling.

Claim For any function g , the expected training error is equal to the generalization error.

QUESTION. Suppose we find g with training error 0. How does this imply that the generalization error is small? Factors that influence this question: Amount of data (how faithful is the sample?). The complexity of the function (# of functions considered). Algorithm used to find g . This is often succinctly phrased as "does g generalize?" If answer is no, this implies that you've "overfit" the data. First, we'll consider the "well-separated finite setting" Here, we will make two enormous assumptions. Assume that: The ground truth labelling function $f \in S = \{f_1, f_2, \dots, f_n\}$. That is, f is a function with a finite number of functions which is finite. All of these functions are well-separated. No function in the class is similar to f . For all $f_i \in S$ with $f_i \neq f$, the generalization error of $f_i > \delta$. Note that this sort of a silly assumption, but we will drop both. Naive "algorithm" return any g in our set S that has training error 0.

Main theorem Given assumptions 1 and 2, if the number of datapoints $n > \frac{1}{\delta^2} (\log k + \log \frac{1}{\delta})$, then, with probability at least $1 - \delta$, g will generalize. Some comments: logarithmic function in k and $1/\delta$ is good, but inverse linear function in $\frac{1}{\delta}$ is kind of bad. We will prove this in two parts. First, we will analyze the probability that we are "tricked" by a bad f_i . Next, we will union bound over all bad f_i s. Consider a bad function f_i . The probability that we are tricked by this function is

$$Pr[\text{TrainingError}(f_i) = 0] = \prod_{j=1}^n Pr(f_i(x_j) = f(x_j))$$
$$\leq (1 - \delta)^n$$
$$< e^{-n\delta}.$$

The least inequality follows from the inequality $1 - x \leq e^{-x}$. Proof from Taylor series / plot. There are at most k "bad" functions. Hence we can apply a union bound, to obtain $\delta = Pr(\text{output bad function}) \leq ke^{-n\delta}$. Now, the desired result directly follows from solving for failure probability $\frac{1}{k}$. Results of this form are generally referred to as being in the "PAC" framework (probably approximately correct).

Theorem for linear classifiers For linear classifiers, if the number of datapoints $n \geq \frac{1}{\delta^2} (d + \log \frac{1}{\delta})$, then, with probability $1 - \delta$, g will generalize.

Generalization Does g generalize is equivalent to asking if its test error is the same as its training error. Start with two assumptions: 1) there are h (finite) possible ground truth functions, and every other candidate has generalization at least ϵ . The probability of being tricked by a single function that matches on all n datapoints is $\leq (1 - \epsilon)^n \leq e^{-\epsilon n}$. Union bounding over h possible functions, get probability $\leq he^{-\epsilon n}$. Letting $n \geq \frac{1}{\epsilon^2} ((nh) + \ln \frac{1}{\delta})$ then with probability at least $1 - \delta$ we get the correct ground truth functions (sample complexity). Even if assumption (2) doesn't hold, can rephrase theorem as $n \geq \frac{1}{\epsilon^2} ((nh) + \ln \frac{1}{\delta})$ implies with probability at least $1 - \delta$ over the samples, the output of the learning algorithm is a function with generalization error less than ϵ (PAC guarantee). Can use curse of dimensionality for arbitrary linear classifiers - only exponentially many directions to look in (exp in dimension). Then if f is a linear classifier, $n \geq \frac{1}{\epsilon^2} (d + \ln \frac{1}{\delta})$, then with probability at least $1 - \delta$ the output of the learning algorithm is a function with generalization error less than ϵ . Rule of thumb: make sure your training data size n is linear in the number d of free parameters/features you're trying to learn. If no function has 0 training error, output the function which has the smallest training error (ERM algorithm). Then if $n \geq \frac{1}{\epsilon^2} (d + \ln \frac{1}{\delta})$, with probability at least $1 - \delta$, the generalization error of the output is the training error $\pm \epsilon$. Functions do not have to be linear classifiers, can be anything. In the same setting, if ϵ is the minimum generalization error of any linear classifier and $n \geq \frac{1}{\epsilon^2} (d + \ln \frac{1}{\delta})$, then the output of the ERM algorithm will have generalization error $\leq \epsilon + 2\epsilon$. Note: from going to 0 percent error to nonzero error (i.e. realizable to non-realizable, we pick up another factor of $\frac{1}{\epsilon}$.

Regularization Kernelization allows you to do polynomial embedding without actually writing out extra terms (do it implicitly instead). If we have a ton of data ($n > d$), then can use polynomial embeddings or random project + nonlinearity (neural networks). If the features in the input space are meaningful (i.e. coordinates matter), probably want polynomial embedding (or if you want to preserve interpretability of features). If your features are rotationally invariant, use random project + nonlinearity. Bayesian viewpoint: true model underlying the data is drawn from a known prior. Assuming gaussian noise, finding max likelihood function is equivalent to L2 regularization. It's unbiased if Bayesian priors will hold - and note that different priors give rise to different regularizers. Frequentists, in contrast, argue that regularization allows you to recover properties of the true model (e.g. sparsity). l_1 regularizer is a great proxy for l_0 regularizer: Given n independent Gaussian vectors from \mathbb{R}^d , consider labels $y_i = c > x_i$ for some vector a with $||a||_0 = s$, then the minimizer of the l_1 regularized objective function will be the vector a with high probability, provided that $n > c \log(d)$ for some constant c . Note the contrast to $O(d)$ data required with no regularization.

3 - REGULARIZATION

Note on last part - is very open ended, at the cutting edge of machine learning research (but don't feel obliged to write pages of analysis).

PUNCLINE from LAST CLASS. If you have a set k of different functions $\{f_1, \dots, f_k\}$, then the "best" one will generalize if $n > O(\log k)$. If we are classifying in d dimensions - we can approximate by set of $\exp(d)$ linear functions. If $n > d$, expect generalization.

REGULARIZATION. A way to express a set of preferences over models. Such a scheme that will take both performance on training data as well as these preferences into account. For example, we can consider L_2 -regularized least squares. Let $x_1, \dots, x_n \in \mathbb{R}^d$, and $y_1, \dots, y_n \in \mathbb{R}$. In this setting, we want to minimize the following objective:

$$\min_a f(a) = \sum_{i=1}^n ((x_i a) - y_i)^2 + \underbrace{\lambda ||a||_2^2}_{\text{regularization term}}.$$

There are two broad types of regularization, explicit or implicit: Explicit regularization (e.g. L_2 regularization, preferring sparse vectors). Implicit regularization (algorithm that has had preferences).

Why should we regularize; why wouldn't we just return the empirical risk minimizer?

PERSPECTIVE. You always want roughly $n \approx d$, where generalization might not hold. If you have $n = 1,000,000$, then you want $d \approx 1,000,000$. Otherwise - it's sort of a "waste"; if it is truly 1000 in your dataset, try to construct additional features to learn a model in 1,000,000 dimensional space.

How should we construct additional features? Here are two approaches.

POLYNOMIAL EMBEDDING. For example - quadratic embedding. $x = (x_1, x_2, \dots, x_d) \rightarrow f(x) = (x_1, \dots, x_d, x_1^2, x_1 x_2, x_1 x_3, \dots, x_1^2, 1) \in \mathbb{R}^{d+1+\binom{d}{2}}$. One simple extension of this would use quadratic features to fit a classifier is when you are fitting a circular decision boundary.

RANDOM PROJECTION + NON-LINEARITY. You really need non-linearity, since otherwise you'll just be learning another "linear" function. Can choose $\sqrt{x}, x^2, \sigma(x)$, or most other "nice" nonlinearities (since they all roughly have similar properties).

Downsides of adding new features: One objection could be that working with $\approx d^2$ dimensional points is annoying. But this isn't an issue: you can "implicitly" work over the embedded points without computing the embedding. The area of math devoted to this is referred to as "kernelization" (usually the context of SVMs).

Real issue: if you need d^2 features, you generally need much more data.

Rule of thumb: if the coordinates actually have significance, the polynomial embedding preserves interpretability. How should you think about regularization? There are two views: the Bayesian view, and the frequentist view.

BAYESIAN VIEW. Assume that the true model is drawn from some known "Prior" distribution. This allows us to evaluate the "likelihood" / probability of a given candidate model.

FREQUENTIST VIEW. Goal: argue that if true model has "nice structure", then can find it.

BAYESIAN APPROACH TO REGULARIZATION. ("Gaussian prior") Suppose we have $x_1, \dots, x_n \in \mathbb{R}^d$, assume that the true label $a \in \mathbb{R}^d$ is drawn by choosing each coordinate independently from $(0, 1)$. Now, suppose that each label is set as $y_i = (x_i, a^T) + z_i$, where noise $z_i \sim (0, \sigma^2)$. Now, given $x_1, \dots, x_n, y_1, \dots, y_n$, ask Likelihood(a) = $Pr(a) Pr(\text{data} | a)$. Because we made strong assumptions on the label distributions, we can directly compute these probabilities. Hence

$$\text{Likelihood}(a) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp(-a_i^2 / 2) \prod_{i=1}^n \exp(-((x_i, a) - y_i)^2 / 2)$$
$$\propto \exp(-||a||_2^2 / 2 - \frac{1}{2\sigma^2} \sum_{i=1}^n ((x_i, a) - y_i)^2)$$

Maximum likelihood of a is equivalent to minimizing $\sum_{i=1}^n ((x_i, a) - y_i)^2 + 2\sigma^2 ||a||_2^2$. This derivation even tells us how to set the regularization constant - should be $2\sigma^2$ times the variance of the noise.

FREQUENTIST APPROACH TO REGULARIZATION / REGULARIZATION. Consider a model a^* that is s -sparse (i.e. there are s nonzero coordinates). Question: why do we care about sparse models? One answer is that lots of models are actually sparse. Think of the laws of physics. Other view: maybe the world is sloppy, and the best model might not be sparse. But what's most helpful for interpretability is fitting a sparse model. Question: can we build a regularizer that lets us selectively find sparse models? The obvious choice is

$f(a) = \sum_{i=1}^n ((x_i, a) - y_i)^2 + \lambda ||a||_0$ where we are using the "0-norm" that computes sparsity.

CLAIM. If $n > O(\log(d))$ then the sparsest model that fits the data is "correct". The number of s -sparse d -dimensional function is just $\binom{d}{s}$, and there are about $\exp(s)$ sparse functions. So there are approximately $d^s \exp(s)$ sparse functions, and we need datapoints around the logarithm of this. The problem with using sparse models is that it is not differentiable (so finding the minimizer is NP-hard). So, we can note the following: l_0 regularization is great, but computationally intractable. Idea: use l_1 regularization as proxy for l_0 . Miraculously, the claim from before still holds for l_1 regularization (proved in the early 2000's, Candès in the stat / math department here).

4 - PCA

Principal component analysis can be used to analyze the structure of a data set or allow the representation of the data in a lower dimensional space(s) as well as many other applications).

Let $\{\tilde{x}_i\}$ be a set of N column vectors of dimension D . Define the scatter matrix S_x of the data set as $S_x = \sum_{i=1}^N ((\tilde{x}_i - \bar{\mu}) - (\tilde{x}_i - \bar{\mu}))^T$ where $\bar{\mu}$ is the mean of the dataset $\bar{\mu} = \frac{1}{N} \sum_{i=1}^N \tilde{x}_i$.

The d largest principal components are the eigenvectors w_i corresponding to the d largest eigenvalues. S can be chosen arbitrarily with $d < D$. The eigenvectors of S can usually be found by using singular value decomposition.

The dominant eigenvectors describe the main directions of variation of the data. For example, if a dataset had 2 large eigenvalues, then the data variation is described largely by the linear combinations of the 2 corresponding eigenvectors (i.e. the data is largely coplanar).

The d eigenvectors can also be used to project the data into a d dimensional space. Define $W = [\tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_d]$. The projection of vector \tilde{x} is $\tilde{y} = W^T \tilde{x}$. The scatter matrix S_y of the vectors $\{\tilde{y}_i\}$ is: $S_y = W^T S_x W$. The matrix W maximizes the determinant of S_y for a given d .

PCA PCA does not necessarily preserve pairwise distances like JL. JL is also data-oblivious, PCA is deterministic. Coordinates used in JL have no intrinsic meaning. PCA often finds latent variables. Linear regression has a "preferred" variable. PCA treats all coordinates equally - works well when there are a set of latent variables and all coordinates are linear combinations of those variables. PCA requires mean subtraction, and scaling of coordinates so that each coordinate has unit standard deviation. PCA assumes that the data corresponds to interesting information. PCA tries to maximize the square projection of each vector onto subspace, alternatively lower its square distances from the hyperplane. PCA does not find nonlinear structure. If the structure is non-orthogonal, PC components may lose interpretability be they're forced to be orthogonal. PCA can also fail due to high-variance noise. Can show with some extra work that we're trying to diagonalize the covariance matrix $X^T X$. Every covariance matrix can be diagonalized by an orthogonal matrix, i.e. $X^T X = Q D Q^T$, Q orthogonal. If only need the top few principal components, compute them with power iteration: choose random unit vector; apply covariance matrix, rescale to unit, continue until convergence. Converged vector will be top principal component. Subtract out this PC from everything and iterate. With probability at least $\frac{1}{2}$ over the choice of initial unit vector, for $t > 1$, m datapoints in n dimensions, $|< A^t a, v_1 >| \geq 1 - 2n^{-1} (\frac{2\lambda}{\lambda_1})^t$, so number of iterations required scales as $\frac{\log m}{\log(\frac{1}{2})}$.

gap. Can exponentially improve this probability $\frac{1}{2}$ by repeating this with different random initial unit vectors.

SVD Low-rank matrix approximations are great for matrix completions. $A = USV^T$. Left singular vectors are in U , right singular vectors are in V (remember diagram). U, V are orthogonal. S is diagonal matrix with entries sorted in nonincreasing order. Only for covariance matrix is $U = V$. The columns of U are called left singular vectors of A and the rows of V^T (columns of V) are the right singular vectors of A). Running time is $\min m^2 n$ or $n^2 m$. If just want k largest singular values and singular vectors, can do it in knn time. Rank- k SVD is optimal in Frobenius norm of all rank- k matrices. How to choose k ? Choose such that top k singular values sum to a multiple of the rest (maybe 10). The right singular vectors of S are the same as the eigenvectors of $X^T X$ (the principal components), and the eigenvalues of the covariance matrix are the square of the singular values of X . So PCA reduces to computing the SVD of A without forming the covariance matrix. If you only need the top few eigenvectors, use Power Iterations If you want them all, use SVD. In Frobenius norm, each row of $A - A_k$ is exactly the distance from each point to the span of the top k principal components. SVD and PCA are equivalent. The rows of X can be interpreted as linear combinations of the rows of V^T . The columns of X can be interpreted as linear combinations of the columns of U (the left singular vectors). When the columns of U are interpretable, this is valuable. Otherwise, use PCA. In customers x products, the right singular vectors are customer types and left singular vectors are product types. Can develop PCA (not SVD) based low-rank approximations by computing covariance, and each row as the projection onto the top k principal components, which is exactly the same as rank- k SVD approximation.

5 - SVD

Rank k matrix rank k if it can be written as the sum of the rank k matrix ones.

Low rank approximation examples Compression, de-noising, matrix-completion.

SVD ASD of an $m \times n$ matrix A can be written as $A = USV^T$, where U is $m \times m$ ortho, V is $n \times n$ ortho, and S is an $m \times n$ diagonal matrix with nonnegative entries (diagonal entries sorted from high to low). Columns of U : left SVs, cols of V : right SVs, and entries of S are singular values. Note that SVD is equivalent to $A = \sum_{i=1}^{\min(m,n)} s_i u_i v_i^T$.

Rank k matrix Rank k if it can be written as the sum of the rank k matrix ones. Rank k approximation is $A_k = U_k S_k V_k^T$. Storing the matrices in rank k is $O(k(m+n))$, compared to $O(mn)$ space for original matrix A . Note that $||A - A_k||_F \leq ||A - B||_F$ for every matrix A , rank target $k \geq 1$, and rank $m \times n$ matrix B .

PCA reduces to SVD We start by sketching the equivalence of truncated rank k SVD and a PCA projection onto the first k principal components. Suppose X is an $n \times p$ matrix (assume it is mean normalized, for simplicity). The $p \times p$ covariance matrix is computed as $X^T X$. Since it is symmetric, we can diagonalize it as $X^T X = Q \Lambda Q^T$ where Q is a matrix of eigenvectors, and Λ is a diagonal matrix with eigenvalues λ_i in decreasing order on the diagonal. Projecting our data onto these eigenvectors yield the principal components.

Now, performing singular value decomposition gives us $X = USV^T$, where U is a unitary matrix, and S is the diagonal matrix of singular values s_i . Expressing the covariance matrix in terms of the singular value decomposition of X , we can write $X^T X = VSU^T USV^T / (n-1) = V \frac{S^2}{n-1} V^T$. Therefore, the right singular vectors V are the eigenvectors of $X^T X$. Furthermore, singular values can be expressed in terms of the eigenvalues of $X^T X$ with the relation $A_i = s_i^2 / (n-1)$.

5B - TENSORS

Tensors Examples: 3-Tensor with works that appear next to each other, moments Tensor, n -Tensor's basis are outer product of n different sets of vectors. FACTS: For tensors, cannot iteratively subtract rank-1 approximations and test rank $k \approx 1$ approximation. Also do not get the rank over the reals is the same as the rank over complex numbers for k -tensors with $k \geq 3$, even if tensor is real-valued. Also don't know how to construct $n \times n \times n$ tensors whose rank is greater than $n^{1/3}$ for all n - hard to reason about rank of tensors. Computing the rank of 3-tensors is NP-hard. But, if the rank of a 3-tensor is sufficiently small, it can be efficiently computed and its low-rank representation is unique and can be efficiently recovered. The factors do NOT need to be orthogonal, as long as they are linearly independent.

6 - SPECTRAL GRAPH THEORY **Graph** Given a graph $G = (V, E)$ with $|V| = n$ vertices, can associate matrices to graph. Laplacian matrix is $L_G = D - A$, where D is degree matrix (diagonal matrix where $D(i, i)$ is degree of i th node in G), and A is the adjacency matrix, with $A(i, j) = 1$ if $(i, j) \in E$. Note that $\sum_j A(i, j) = \sum_{(i, j) \in E} (v(i) - v(j))^2$, v assigns number of each vertex in G , and $\sum_j A(i, j) =$ sum of squares of differences between values of neighboring nodes. **Eigenspaces** From the previous identity, eigenvalues are nonnegative. Number of zero eigenvalues of L is the number of connecting components of G . The maximum eigenvalue $\max_{i=1 \dots n} \lambda_i$ will maximize the discrepancy between neighbors values of v .

Spectral embeddings Smallest eigenvector is not helpful, since $v_1 = (\frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}})$. Helpful to consider v_2 and v_3 , and plot these together. **Isoperimetric ratio** Given a set S , the isoperimetric ratio $\theta(S)$ is defined as $\theta(S) = \frac{|\partial(S)|}{\min_{i=1 \dots n} \lambda_i \theta(S)}$. The isoperimetric number is $\theta_n = \min_{i=1 \dots n} \lambda_i \theta(S)$. **Conductance** The conductance of a partition of a graph is defined as $\text{cond}(S) = \frac{|\partial(S)|}{\min_{i=1 \dots n} \sum_{j \in S} \deg(j) \lambda_i \theta(S)}$.

λ_2 Isoperimetric number of graph satisfies $\theta_n \geq \lambda_2 (1 - \frac{1}{|\partial(S)|})$.

Partial converse (Cheeger) If λ_2 is second smallest eigenvalue, there exists a set $S \subset V$ such that $\text{cond}(S) \leq \frac{\sqrt{2\lambda_2}}{\sqrt{d}}$.

7 - SAMPLING / ESTIMATION

R - problem Want to sample k uniformly random elements from a datastream of length $N \gg k$. N large, and N unknown. **Reservoir sampling** Given a number k , and a datastream x_1, x_2, \dots of length $> k$: Put the first k elements of the stream into a reservoir $R = \{x_1, \dots, x_k\}$. For $i \geq k + 1$, with probability $\frac{k}{i}$ replace a random entry of R with x_i . Return reservoir at the end of stream.

RS is uniform For $i \geq k$, $Pr[x_i \in R_k] = \frac{k}{i}$ where R_k is the reservoir after time i . Proof by induction.

Markov Let X be a real-valued RV with $X \geq 0$. For any $c > 0$, $Pr[X \geq cE[X]] \leq \frac{1}{c}$. Intuition: at most 10% of population can have income $> 10 \times$ avg. income.

Chebyshev Let X be real-valued RV, and $c > 0$. Then $Pr[|X - E[X]| \geq c \sqrt{\text{Var}[X]}] \leq \frac{1}{c^2}$.

Importance sampling We can estimate properties of distribution A , based on samples from distribution B . Idea: have distribution B place higher weight on the "important" elements of the domain.

Importance sampling ex. For example, suppose we know that computer scientists have higher income, and that computer scientists are 0.05 of population. Instead of taking n random samples from pop, take 0.8n samples of non-computer scientists, and 0.2 samples of computer scientists. Then: $\text{avg} = 0.95a_1 + 0.05a_2$ allows us to estimate the population average salary. Benefit: reduce variance of the estimate, by focussing our samples on the important portion of the distribution.

Good-Turing frequency estimation Given $n</$

Theorem comments Extends to “almost k -sparse” signals. Extends to non-Gaussian matrices A . However, sparse matrices d don't work. m is optimal up to a constant factor, since $\log_2 \binom{n}{k} = O(k \log \frac{n}{k})$ by Stirling's approximation.

Uncertainty principle No signal can be made sparse simultaneously in both time and frequency domains. **LPs** Consist of decision variables, linear constraints of the form $\sum_{j=1}^n a_{ij}x_j \leq b_i$, or $\sum_{j=1}^n a_{ij}x_j = b_i$. Objective function should be linear, of the form $\min \sum_{j=1}^n c_jx_j$.

ℓ_1 **linearization** Add y_j to represent $|x_j|$. Use objective function $\min \sum_{j=1}^n y_j$ which is linear, with $2n$ inequalities of the form $y_j - x_j \leq 0, y_j + x_j \geq 0$. Can also extend LPs to accommodate noise.

Convexity A set $C \subseteq \mathbb{R}^n$ is convex for every $x, y \in C$ and $\lambda \in [0, 1], \lambda x + (1 - \lambda)y \in C$. A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if the region above f is a convex set. That is, $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$ for every $x, y \in \mathbb{R}^n$ and $\lambda \in [0, 1]$. Note that LPs are convex, since intersection of half-spaces will be convex. Also, the set of $n \times n$ symmetric and PSD matrices viewed as a subset of \mathbb{R}^{n^2} is convex.

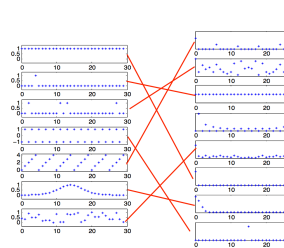
Matrix completion Suppose we have a matrix \hat{M} derived from M by erasing its entries. Want to recover M , with key assumption that M has low rank. Suppose M has SVD $U\Sigma V^T$. Rank of M is the number of non-zero singular values, with every row of M a linear combination of top r right SVs, and every col is a linear combination of top r left SVs. Optimization problem can be viewed as $\min \|\text{supp}(\Sigma(M))\|_1$, subject to M agrees with \hat{M} on known entries. The exact problem is ℓ_0 , but we can relax it to optimization the ℓ_1 norm with $\min \|\Sigma(M)\|_1$ (which is convex).

Matrix completion theorem If M has rank r , and \hat{M} has at least $\Omega(r(m+n)\log^2(m+n))$ known entries, chosen uniformly at random from M . Then M is sufficiently dense and non-pathological.

Application: Matrix Completion Key assumption: we want to recover the entire $n \times n$ partially observed matrix M (observed is M_{hat}) that **HAS LOW RANK**. If we try to minimize the rank of M subject to M agreeing with M_{hat} on observed entries, we get NP-hard problem. But if we reframe, we can say we want to minimize the number of nonzero singular values of M s.t. M agrees with M_{hat} on known entries. Can again change this to **nuclear norm minimization** where we want to minimize the sum of the singular values and solve: $\min \|\Sigma(M)\|_1$. s.t. M agrees with M_{hat} on known entries. In fact, the objective function here is convex and can be solved efficiently.

ASSIGNMENTS

Let X, Y be discrete integer random variables. Then we can write the distribution of $Z = X + Y$ as a convolution between X and Y : $P(Z = z) = \sum_{x=-\infty}^{\infty} P(X = k)P(Y = z - k)$. In particular, the vector $\mathbf{q} = p \circ s \circ p \circ \dots \circ p$ (convolved 100 times) represents the probability distribution over 100 rolls of a six sided die. Since $\mathbf{q}[0]$ represents the probability that we roll a total of 100, $\mathbf{q}[150]$ will represent the probability that the 100 rolls sum to 250.



Let \mathbf{f}, \mathbf{g} be N -tuples. Let $\mathbf{h} = \mathbf{f} * \mathbf{g}$. Computing the Fourier transform, we obtain

$$\begin{aligned} \mathbf{F}_h(m) &= \sum_{j=0}^{2N-1} h[j] \exp(-2\pi i m j / N) \\ &= \sum_{j=0}^{2N-1} \left(\sum_{k=0}^{2N-1} f[k] g[j-k] \right) \exp(-2\pi i m j / N) \\ &= \sum_{k=0}^{2N-1} f[k] \left(\sum_{j=0}^{2N-1} g[j-k] \exp(-2\pi i m j / N) \right) \\ &= \sum_{k=0}^{2N-1} f[k] \left(\exp(-2\pi i m k / N) \mathcal{F}(g^*) \right) \\ &= \mathcal{F}(g^*) \left(\sum_{k=0}^{2N-1} f[k] \exp(-2\pi i m k / N) \right) \\ &= \mathcal{F}(f^*) \mathcal{F}(g^*), \end{aligned}$$

which proves the desired result. Therefore, we can implement the convolution $\mathbf{f} * \mathbf{g}$ in the frequency domain as follows: $\mathbf{f} * \mathbf{g} = \mathbf{F}^{-1} (\mathbf{F}^* \cdot \mathbf{g}^*)$. A naive implementation of convolution would require $O(N^2)$ operations. Since the complexity of the FFT is $O(n \log n)$, a FFT-based convolution is also $O(n \log n)$, and therefore much more efficient.

If \mathbf{f} is an M -tuple and \mathbf{g} is an M -tuple with $M < N$ we can zero-pad \mathbf{f} and \mathbf{g} to length $M + N$, and perform the same computation.

PCA does poorly with noise only in Y because it aims to capture the variance in the y direction because of noise. As the noise coefficient c increases, PCA will upweight the principal component in the y direction, which increases the slope as c increases.

Formally, the first principal component satisfies

$$\begin{aligned} \mathbf{w}_1 &= \|\mathbf{w}\| = 1 \left\{ \sum_i (\mathbf{x}_i \cdot \mathbf{w})^2 \right\} \\ &= \|\mathbf{w}\| = 1 \left\{ \sum_i \left((i/1000, 2i/1000 + \mathcal{N}(0, c)) \cdot \mathbf{w} \right)^2 \right\}. \end{aligned}$$

Since the noise is biased towards the second component of x_i , the net effect is to cause each dot product, biasing the first principal component towards the y direction.

PCA does well with noise in both X and Y , because there is now variance in both dimensions. In this setting, we can compute \mathbf{w}_1 as follows:

$$\begin{aligned} \mathbf{w}_1 &= \|\mathbf{w}\| = 1 \left\{ \sum_i (\mathbf{x}_i \cdot \mathbf{w})^2 \right\} \\ &= \|\mathbf{w}\| = 1 \left\{ \sum_i \left((i/1000 + \mathcal{N}(0, c), 2i/1000 + \mathcal{N}(0, c)) \cdot \mathbf{w} \right)^2 \right\}. \end{aligned}$$

Since the variance is spread across both dimensions, \mathbf{w} does not get “biased” towards either dimension, and PCA recovery results in a good estimate of the slope.

Least squares performs poorly because with increasing noise in X and Y , there is less of a predictive relationship between X and Y .

The least squares slope estimate can be computed as follows:

$$\begin{aligned} (X_n - \bar{X}_n, Y_n - \bar{Y}_n) / \|(X_n - \bar{X}_n)\|_2^2, \\ \text{where} \\ X_n = X + \mathcal{N}(0, c) \\ Y_n = Y + \mathcal{N}(0, c). \end{aligned}$$

This results in a decreasing slope as the noise coefficient c increases. In particular, as $c \rightarrow \infty$, the slope $\rightarrow 0$, since there will be no correlation between X and Y .

NEW STUFF

Frequentists vs. Bayesians what is probability? One is called the **frequentist** interpretation. In this view, probabilities represent long run frequencies of events. For example, the above statement means that, if we flip the coin many times, we expect it to land heads about half the time. The other interpretation is called the **Bayesian** interpretation of probability. In this view, probability is used to quantify our **uncertainty** about something; hence it is fundamentally related to information rather than repeated trials (Jaynes 2003). In the Bayesian view, the above statement means we believe the coin is equally likely to land heads or tails on the next toss. The big advantage of the Bayesian interpretation is that it can be used to model our uncertainty about events that do not have long term frequencies. For example, we might want to compute the probability that the polar ice cap will melt by 2020 CE. This event will happen zero or one times, but cannot happen repeatedly. Nevertheless, we ought to be able to quantify our uncertainty about this event. To give another machine learning oriented example, we want to compute a “blip” on our radar screen, and want to compute the probability distribution over the location of the corresponding target (be it a bird, plane, or missile). In all these cases, the idea of repeated trials does not make sense, but the Bayesian interpretation is valid and indeed quite natural. We shall therefore adopt the Bayesian interpretation in this book. Fortunately, the bare rules of probability theory are the same, no matter which interpretation is adopted.

A brief review of probability theory
Basic concepts We denote a random event by defining a **random variable** X .
Discrete random variable: X can take on any value from a finite or countably infinite set.
Continuous random variable: the value of X is real-valued.
CDF

$$F(x) \triangleq P(X \leq x) = \begin{cases} \sum_{x_i \leq x} p(u) & , \text{discrete} \\ \int_{-\infty}^x f(u) du & , \text{continuous} \end{cases} \quad (1)$$

PMF and PDF For discrete random variable, We denote the probability of the event that $X = x$ by $P(X = x)$, or just $p(x)$ for short. Here $p(x)$ is called a **probability mass function** or **PMF**. A probability mass function is a function that gives the probability that a discrete random variable is exactly equal to some value¹. This satisfies the properties $0 \leq p(x) \leq 1$ and $\sum_{x \in \mathcal{X}} p(x) = 1$. For continuous variables, in the equation $F(x) = \int_{-\infty}^x f(u) du$, the function $f(x)$ is called a **probability density function** or **PDF**. A probability density function is a function that describes the relative likelihood for this random variable to take on a given value². This satisfies the properties $f(x) \geq 0$ and $\int_{-\infty}^{\infty} f(x) dx = 1$.

Multivariate random variables
Joint CDF We denote joint CDF by $F(x, y) \triangleq P(X \leq x \cap Y \leq y) = P(X \leq x, Y \leq y)$.

$$F(x, y) \triangleq P(X \leq x, Y \leq y) = \sum_{x' \leq x} \sum_{y' \leq y} p(u, v) = \int_{-\infty}^x \int_{-\infty}^y f(u, v) du dv \quad (2)$$

$$\text{product rule:} \quad p(X, Y) = P(X|Y)P(Y) \quad (3)$$

$$\begin{aligned} \text{Chain rule:} \\ p(X_{1:N}) &= p(X_1)p(X_2|X_1)p(X_3|X_2, X_1) \dots p(X_N|X_{1:N-1}) \end{aligned} \quad (4)$$

$$\begin{aligned} \text{Marginal distribution Marginal CDF:} \\ F_X(x) &\triangleq F(x, +\infty) = \\ &= \sum_{x' \leq x} P(X = x_i) = \sum_{x' \leq x} \sum_{y' \leq +\infty} P(X = x_i, Y = y_j) \\ &= \int_{-\infty}^x \int_{-\infty}^{+\infty} f(u, v) du dv \end{aligned} \quad (5)$$

$$\begin{aligned} F_Y(y) &\triangleq F(+\infty, y) = \\ &= \sum_{x' \leq +\infty} p(Y = y_j) = \sum_{x' \leq +\infty} \sum_{y' \leq y} P(X = x_i, Y = y_j) \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^y f(u, v) dv du \end{aligned} \quad (6)$$

$$\begin{aligned} \text{Marginal PMF and PDF:} \\ \begin{cases} P(X = x_i) = \sum_{j=1}^m P(X = x_i, Y = y_j) & , \text{ discrete} \\ f_X(x) = \int_{-\infty}^{+\infty} f(x, y) dy & , \text{ continuous} \end{cases} \end{aligned} \quad (7)$$

$$\begin{cases} p(Y = y_j) = \sum_{i=1}^m P(X = x_i, Y = y_j) & , \text{ discrete} \\ f_Y(y) = \int_{-\infty}^{+\infty} f(x, y) dx & , \text{ continuous} \end{cases} \quad (8)$$

$$\text{Conditional distribution Conditional PMF:} \quad p(X = x_i | Y = y_j) = \frac{p(X = x_i, Y = y_j)}{p(Y = y_j)} \text{ if } p(Y) > 0 \quad (9)$$

The pmf $p(X|Y)$ is called **conditional probability**.
Conditional PDF:

$$f_{X|Y}(x|y) = \frac{f(x, y)}{f_Y(y)} \quad (10)$$

$$\begin{aligned} \text{Bayes rule} \\ p(Y = y | X = x) &= \frac{p(X = x, Y = y)}{p(X = x)} \\ &= \frac{p(X = x | Y = y) p(Y = y)}{\sum_{y'} p(X = x | Y = y') p(Y = y')} \end{aligned} \quad (11)$$

Independence and conditional independence We say X and Y are unconditionally independent or marginally independent, denoted $X \perp Y$, if we can represent the joint as the product of the two marginals, i.e.,

$$\begin{aligned} X \perp Y &= P(X, Y) = P(X)P(Y) \quad (12) \\ \text{We say } X \text{ and } Y \text{ are conditionally independent(CI) given } Z \text{ if} \\ \text{the conditional joint can be written as a product of conditional marginals:} \\ X \perp Y | Z &= P(X, Y | Z) = P(X | Z)P(Y | Z) \quad (13) \end{aligned}$$

Quantiles Since the cdf F is a monotonically increasing function, it has an inverse; let us denote this by F^{-1} . If F is the cdf of X , then $F^{-1}(\alpha)$ is the value of x_α such that $P(X \leq x_\alpha) = \alpha$; this is called the α quantile of F . The value $F^{-1}(0.5)$ is the **median** of the distribution, with half of the probability mass on the left, and half on the right. The values $F^{-1}(0.25)$ and $F^{-1}(0.75)$ are the lower and upper **quartiles**.

Mean and variance The most familiar property of a distribution is its **mean**, or **expected value**, denoted by μ . For discrete rv's, it is defined as $\mathbb{E}[X] \triangleq \sum_{x \in \mathcal{X}} xp(x)$, and for continuous rv's, it is defined as $\mathbb{E}[X] \triangleq \int_{\mathcal{X}} xp(x)dx$. If this integral is not finite, the mean is not defined (we will see some examples of this later). The **variance** is a measure of the “spread” of a distribution, denoted by σ^2 . This is defined as follows:

$$\begin{aligned} \text{var}[X] &= \mathbb{E}[(X - \mu)^2] \\ &= \int (x - \mu)^2 p(x) dx \\ &= \int x^2 p(x) dx + \mu^2 \int p(x) dx - 2\mu \int xp(x) dx \\ &= \mathbb{E}[X^2] - \mu^2 \end{aligned} \quad (14) \quad (15)$$

$$\begin{aligned} \text{from which we derive the useful result} \\ \mathbb{E}[X^2] &= \sigma^2 + \mu^2 \end{aligned} \quad (16)$$

$$\begin{aligned} \text{The standard deviation is defined as} \\ \text{std}[X] &\triangleq \sqrt{\text{var}[X]} \end{aligned} \quad (17)$$

This is useful since it has the same units as X itself.
Some common discrete distributions In this section, we review some commonly used parametric distributions defined on discrete state spaces, both finite and countably infinite.
The Bernoulli and binomial distributions Now suppose we toss a coin any number of times. Let $X \in \{0, 1\}$ be a binary random variable, with probability of “success” or “heads” of θ . We say that X has a **Bernoulli distribution**. This is written as $X \sim \text{Ber}(\theta)$, where the pmf is defined as

$$\text{Ber}(x|\theta) \triangleq \theta^{\mathbb{I}(x=1)} (1 - \theta)^{\mathbb{I}(x=0)} \quad (18)$$

Suppose we toss a coin n times. Let $X \in \{0, 1, \dots, n\}$ be the number of heads. If the probability of heads is θ , then we say X has a **binomial distribution**, written as $X \sim \text{Bin}(n, \theta)$. The pmf is given by

$$\text{Bin}(k|n, \theta) \triangleq \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad (19)$$

The multinoulli and multinomial distributions The Bernoulli distribution can be used to model the outcome of one coin tosses. To model the outcome of tossing a K -sided die, let $\vec{x} = (\mathbb{I}(x = 1), \dots, \mathbb{I}(x = K)) \in \{0, 1\}^K$ be a random vector (this is called **multinoulli encoding** or **one-hot encoding**), then we say X has a **multinoulli distribution** (or categorical distribution³), written as $X \sim \text{Cat}(\theta)$. The pmf is given by:

$$p(\vec{x}) \triangleq \prod_{k=1}^K \theta_k^{\mathbb{I}(x=k)} \quad (20)$$

Suppose we toss a K -sided die n times. Let $\vec{x} = (x_1, x_2, \dots, x_K) \in \{0, 1, \dots, n\}^K$ be a random vector, where x_j is the number of times side j of the dice occurs, then we say X has a **multinomial distribution**, written as $X \sim \text{Mu}(n, \vec{\theta})$. The pmf is given by

$$p(\vec{x}) \triangleq \binom{n}{x_1 \dots x_K} \prod_{k=1}^K \theta_k^{x_k} \quad (21)$$

where $\binom{n}{x_1 \dots x_K} \triangleq \frac{n!}{x_1! x_2! \dots x_K!}$. Bernoulli distribution is just a special case of a Binomial distribution with $n = 1$, and so is multinoulli distribution as to multinomial distribution. See Table ?? for a summary.

The Poisson distribution We say that $X \in \{0, 1, 2, \dots\}$ has a **Poisson distribution** with parameter $\lambda > 0$, written as $X \sim \text{Poi}(\lambda)$, if its pmf is

$$p(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!} \quad (22)$$

The first term is just the normalization constant, required to ensure the distribution sums to 1. The Poisson distribution is often used as a model for counts of rare events like radioactive decay and traffic accidents.

The empirical distribution The **empirical distribution function**⁴, or **empirical cdf**, is the cumulative distribution function associated with the empirical measure of the sample. Let $\mathcal{S} = \{x_1, x_2, \dots, x_N\}$ be a sample set, it is defined as

$$F_N(x) \triangleq \frac{1}{N} \sum_{i=1}^N \mathbb{I}(x_i \leq x) \quad (23)$$

Some common continuous distributions In this section we present some commonly used univariate (one-dimensional) continuous probability distributions.

Gaussian (normal) distribution If $X \sim \mathcal{N}(0, 1)$, we say X follows a **standard normal distribution**. The Gaussian distribution is the most widely used distribution in statistics. There are several reasons for this.

- First, it has two parameters which are easy to interpret, and which capture some of the most basic properties of a distribution, namely its mean and variance.
- Second, the central limit theorem (Section TODO) tells that sums of independent random variables have an approximately Gaussian distribution, making it a good choice for modeling residual errors or “noise”.
- Third, the Gaussian distribution makes the least number of assumptions (has maximum entropy), subject to the constraint of having a specified mean and variance, as we show in Section TODO; this makes it a good default choice in many cases.
- Finally, it has a simple mathematical form, which results in easy to implement, but often highly effective, methods, as we will see.

See (Jaynes 2003, ch 7) for a more extensive discussion of why Gaussians are so widely used.

If $\mathbf{v} = 1$, this distribution is known as the **Cauchy** or **Lorentz** distribution. This is notable for having such heavy tails that the integral that defines the mean does not converge.

To ensure finite variance, we require $\mathbf{v} > 2$. It is common to use $\mathbf{v} = 4$, which gives good performance in a range of problems (Lange et al. 1989). For $\mathbf{v} \gg 5$, the Student distribution rapidly approaches a Gaussian distribution and loses its robustness properties. Here μ is a location parameter and $b > 0$ is a scale parameter. See Figure ?? for a plot. Its robustness to outliers is illustrated in Figure ?? . It also put mores probability density at 0 than the Gaussian. This property is a useful way to encourage sparsity in a model, as we will see in Section TODO.

Introduction Linear regression is the “work horse” of statistics and (supervised) machine learning. When augmented with kernels or other forms of basis function expansion, it can model also nonlinear relationships. And when the Gaussian output is replaced with a Bernoulli or multinoulli distribution, it can be used for classification, as we will see below. So it pays to study this model in detail.

$$\begin{aligned} \text{Representation} \\ p(y|\vec{x}, \vec{\theta}) &= \mathcal{N}(y|\vec{y}^T \vec{x}, \sigma^2) \end{aligned} \quad (24)$$

where \vec{w} and \vec{x} are extended vectors, $\vec{x} = (1, x)$, $\vec{w} = (b, \mathbf{w})$. Linear regression can be made to model non-linear relationships by replacing \vec{x} with some non-linear function of the inputs, $\phi(\vec{x})$

$$p(y|\vec{x}, \vec{\theta}) = \mathcal{N}(y|\vec{y}^T \phi(\vec{x}), \sigma^2) \quad (25)$$

This is known as **basis function expansion**. (Note that the model is still linear in the parameters \vec{w} , so it is still called linear regression; the importance of this will become clear below.) A simple example are polynomial basis functions, where the model has the form

$$\phi(x) = (1, x, \dots, x^d) \quad (26)$$

MLE Instead of maximizing the log-likelihood, we can equivalently minimize the **negative log likelihood** or **NLL**:

$$\text{NLL}(\vec{\theta}) \triangleq -\ell(\vec{\theta}) = -\log(\mathcal{L}(\vec{\theta})) \quad (27)$$

The NLL formulation is sometimes more convenient, since many optimization software packages are designed to find the minima of functions, rather than maxima. Now let us apply the method of MLE to the linear regression setting. Inserting the definition of the Gaussian into the above, we find that the log likelihood is given by

$$\begin{aligned} \ell(\vec{\theta}) &= \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma}} \exp \left(-\frac{1}{2\sigma^2} (y_i - \vec{w}^T \vec{x}_i)^2 \right) \right] \\ &= -\frac{1}{2\sigma^2} \text{RSS}(\vec{w}) - \frac{N}{2} \log(2\pi\sigma^2) \end{aligned} \quad (28) \quad (29)$$

$$\begin{aligned} \text{RSS stands for residual sum of squares and is defined by} \\ \text{RSS}(\vec{w}) &\triangleq \sum_{i=1}^N (y_i - \vec{w}^T \vec{x}_i)^2 \end{aligned} \quad (30)$$

We see that the MLE for \vec{w} is the one that minimizes the RSS, so this method is known as **least squares**. Let's do constants wrt \vec{w} and NLL can be written as

$$\text{NLL}(\vec{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \vec{w}^T \vec{x}_i)^2 \quad (31)$$

There two ways to minimize $\text{NLL}(\vec{w})$.

$$\text{OLS Defining } \vec{y} = (y_1, y_2, \dots, y_N), \vec{X} = \begin{pmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_N^T \end{pmatrix}, \text{ then}$$

$\text{NLL}(\vec{w})$ can be written as

$$\text{NLL}(\vec{w}) = \frac{1}{2} (\vec{y} - \vec{X}\vec{w})^T (\vec{y} - \vec{X}\vec{w}) \quad (32)$$

When \mathcal{S} is small (for example, $N < 1000$), we can use the following equation to compute directly

$$\hat{\vec{w}}_{\text{OLS}} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y} \quad (33)$$

The corresponding solution $\hat{\vec{w}}_{\text{OLS}}$ to this linear system of equations is called the **ordinary least squares** or **OLS** solution. We now state without proof some facts of matrix derivatives (we won't need all of these at this section).

$$\begin{aligned} \text{Tr } A &\triangleq \sum_{i=1}^n A_{ii} \\ \frac{\partial}{\partial A} AB &= B^T \\ \frac{\partial}{\partial A^T} f(A) &= \left[\frac{\partial}{\partial A} f(A) \right]^T \end{aligned} \quad (34) \quad (35)$$

$$\frac{\partial}{\partial A} ABA^T C = CAB + C^T A B^T \quad (36)$$

$$\frac{\partial}{\partial A} |A| = |A| (A^{-1})^T \quad (37)$$

$$\begin{aligned} \text{Then,} \\ \text{NLL}(\vec{w}) &= \frac{1}{2N} (\vec{X}\vec{w} - \vec{y})^T (\vec{X}\vec{w} - \vec{y}) \\ \frac{\partial \text{NLL}}{\partial \vec{w}} &= \frac{1}{2} \frac{\partial}{\partial \vec{w}} (\vec{w}^T \vec{X}^T \vec{X} \vec{w} - \vec{w}^T \vec{X}^T \vec{y} - \vec{y}^T \vec{X} \vec{w} + \vec{y}^T \vec{y}) \\ &= \frac{1}{2} \frac{\partial}{\partial \vec{w}} (\vec{w}^T \vec{X}^T \vec{X} \vec{w} - \vec{w}^T \vec{X}^T \vec{y} - \vec{y}^T \vec{X} \vec{w}) \\ &= \frac{1}{2} \text{Tr} (\vec{w}^T \vec{X}^T \vec{X} \vec{w} - \vec{w}^T \vec{X}^T \vec{y} - \vec{y}^T \vec{X} \vec{w}) \\ &= \frac{1}{2} \frac{\partial}{\partial \vec{w}} (\text{Tr} \vec{w}^T \vec{X}^T \vec{X} \vec{w} - 2\text{Tr} \vec{y}^T \vec{X} \vec{w}) \end{aligned}$$

Combining Equations 36 and 37, we find that

$$\begin{aligned} \frac{\partial}{\partial A^T} ABA^T C &= B^T A^T C^T + BA^T C \\ \text{Let } A^T &= \vec{w}, B = B^T = \vec{X}^T \vec{X}, \text{ and } C = I, \text{ Hence,} \end{aligned}$$

$$\begin{aligned} \frac{\partial \text{NLL}}{\partial \vec{w}} &= \frac{1}{2} (\vec{X}^T \vec{X} \$$