

Machine Learning 1

1 Linear Regression

- Supervised Learning (regression, classification)
- Unsupervised Learning (clustering, dimensionality reduction)
- Matrix: single samples are rows
- Derivative of vector input function is column vector
- $\nabla_x \mathbf{A}x = \mathbf{A}^\top$, $\nabla_x x^\top x = 2x$, $\nabla_x x^\top \mathbf{A}x = 2\mathbf{A}x$
- Linear Regression: fit line $y = f(x) + \varepsilon = w_0 + w_1x + \varepsilon$ (Gaussian noise $\varepsilon \sim N(0, 1)$)
- Minimize summed/mean squared error $\text{SSE} = \sum_{i=1}^N (y_i - f(x_i))^2$ (differentiable, easy to optimize, estimates mean of target function)
- Multiple inputs: $\text{SSE} = (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$ with $\mathbf{X} = \begin{bmatrix} 1 & \mathbf{x}_1^\top \\ \vdots & \vdots \\ 1 & \mathbf{x}_n^\top \end{bmatrix}$
- Least squares solution ($\nabla_{\mathbf{w}} \text{SSE} = 0$): $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ - closed form because SSE convex for linear $f(x)$ (one minimum) and quadratic in w (easy to obtain)
- $R^2 = 1 - \frac{\sum (y_n - \hat{y}_n)^2}{\sum (y_n - \bar{y})^2}$ (quality: how much variation in y explained by variation in x)
- Generalized: $f(x) = \tilde{\mathbf{x}}^\top \mathbf{w} \rightarrow f(x) = \phi(x)^\top \mathbf{w}$, still linear in w (ϕ_i : basis functions)
- $\mathbf{w}^* = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$ with $\Phi = \begin{bmatrix} \phi_1^\top \\ \vdots \\ \phi_n^\top \end{bmatrix}$ (learn any function with suitable ϕ_i)
- Overfitting: model too complex, fits noise / Underfitting: model too simple for data
- Regularization (limit model): Regularization term in cost function with factor λ
- $L_{\text{ridge}} = (\mathbf{y} - \Phi\mathbf{w})^\top (\mathbf{y} - \Phi\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}$ (weight decay / ridge regression)
- $\mathbf{w}_{\text{ridge}}^* = (\Phi^\top \Phi + \lambda \mathbf{I})^{-1} \Phi^\top \mathbf{y}$ (easier to invert due to full rank)

2 Linear Classification

- Expectation of function wrt a distribution: $\mathbb{E}_p[f(x)] = \int p(x)f(x)dx$
- Conditional expectation: $\mathbb{E}_p[f(x)|Y = y] = \int p(x|y)f(x)dx$
- Chain rule: $\mathbb{E}_p[f(x)] = \int p(y)\mathbb{E}_p[f(x)|Y = y]dy$
- Monte-carlo: estimate expectation by samples
- Covariance: $\Sigma = \mathbb{E}_p[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]$, diagonal: variability, other: correlation
- Bernoulli distribution: $p(x) = \mu^x(1 - \mu)^{(1-x)}$ (coin toss)
- Multinomial / Categorical Distribution: $p(c) = \prod \mu_k^{h_{c,k}}$ with 1-hot-encoding (die)
- Gaussian Distribution: $p(x) = N(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-\frac{(x-\mu)^2}{2\sigma^2}\}$
- Multivariate: $p(\mathbf{x}) = N(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\{-\frac{(\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}{2}\}$
- Maximum Likelihood Estimation: $\boldsymbol{\theta}_{\text{ML}} = \text{argmax}_{\boldsymbol{\theta}} \text{loglik}(\boldsymbol{\theta}, D)$
- Linear Gaussian model $p_{\boldsymbol{\theta}}(y|\mathbf{x}) = N(y|\mathbf{w}^\top \tilde{\mathbf{x}}, \sigma^2) \rightarrow$ MLE solution equivalent to least squares, but variance can also be obtained
- Generative model: Assume form of $p(c)$, $p(x|c)$, learn them, predict: compute $p(c|x) \rightarrow$ learn full joint distribution of data (hard), gaussian assumption \rightarrow error
- Discriminative Model: Assume form of $p(c|x)$ and estimate parameters directly from data (simpler than generative modelling, only considers points on border)
- Linear classifier: $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ (\mathbf{w} normal to line, b is bias)
- Counting number of misclassifications as loss is very difficult to optimize (NP-hard)
- Regression loss is not robust to outliers (labels restricted to 0, 1)
- Solution: squash output with sigmoid (bounded between 0 and 1)
- Probabilistic View: $p(c|\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b)^c (1 - \sigma(\mathbf{w}^\top \mathbf{x} + b))^{1-c} \rightarrow$ optimize loglikelihood (logistic regression): cross-entropy loss: $-\sum_i c_i \log f(\mathbf{x}_i) + (1 - c_i) \log(1 - f(\mathbf{x}_i)) \rightarrow$ convex but no closed form \rightarrow gradient descent
- Generalized: Use basis functions to make data linear separable in feature space
- L2 regularization loss: $\text{penalty}(\tilde{\mathbf{w}}) = \|\tilde{\mathbf{w}}\|^2$
- General optimization form: $\text{argmin}_{\boldsymbol{\theta}} \sum l(\mathbf{x}_i, \boldsymbol{\theta}) + \lambda \text{penalty}(\boldsymbol{\theta})$
- Gradient descent: $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$ with learning rate η

- Terminate: change small, gradient small, change in value small, or after fixed time
- Stochastic: use one sample for step (good far away, struggle to find exact optimum)
- Stochastic approximation theory: SGD converges to optimum for strictly convex functions if $\sum \eta_t = \infty$ and $\sum \eta_t^2 < \infty$ (for example $\eta_t = \frac{1}{t}$)
- Stochastic gradients often better than batch since data-set contains redundancy
- Mini-Batches: intermediate between stochastic and batch, preferable for GPU
- Softmax: $p(c = i|\mathbf{x}) = \frac{\exp(\mathbf{w}_i^\top \phi(\mathbf{x}))}{\sum_k \exp(\mathbf{w}_k^\top \phi(\mathbf{x}))}$ (each class gets a weight vector)
- Multiclass Classification: $p(c|\mathbf{x}) = \Pi p(c = k|\mathbf{x})^{\mathbf{h}_{c,k}}$ (use softmax) \rightarrow can be optimized by gradient descent

3 Model Selection

- Model complexity (linear regression: features, λ ; decision tree: depth, number of leaves; NNs: layers, neurons; SVMs: features, regularization; Gaussian Processes: kernel bandwidth) \rightarrow model selection problem
- True risk (unknown) vs. empirical risk (can be evaluated)
- Overfitting (small empirical/high true risk)/Underfitting (high empirical/true risk)
- Expected loss for model and data set size n : $R(\hat{f}_{D_n}) = \mathbb{E}_{D_n}[\mathbb{E}_{x,y}[(\hat{f}_{D_n}(\mathbf{x}) - y)^2]]$
 $= \mathbb{E}_{D_n}[\mathbb{E}_x[(\hat{f}_{D_n}(\mathbf{x}) - \hat{f}_*(\mathbf{x}))^2]] + \mathbb{E}_x[(\hat{f}_*(\mathbf{x}) - f(\mathbf{x}))^2] + \sigma^2$
 $= \text{Variance} + \text{Bias}^2 + \text{Noise}$ (bias (structure error) due to restriction of model, variance due to randomness of data set, \hat{f}_{D_n} : estimate of f from data D_n , $\hat{f}_*(\mathbf{x}) = \mathbb{E}_{D_n}[\hat{f}_{D_n}(\mathbf{x})]$: best model possible, $y = f(\mathbf{x}) + \varepsilon$)
- Underfitting: low variance / high bias; Overfitting: high variance / low bias
- Hold-out method: Judge generalization error with validation data set to pick model (needs more data, unlucky splits can give misleading results)
- Cross validation: Split dataset into k folds, use each as validation once
- Leave-One-Out: $k = n$ / Random sub-sampling: Random points used in each fold
- Avoid overfitting: low complexity, regularize, early stopping, noise, augmentation
- Regularization: penalty $L_2(\theta) = \|\theta\|_2$ (optimizing easy)/penalty $L_1(\theta) = \|\theta\|_1$ (hard, leads to sparse solutions)
- Early stopping: similar effects to L2, efficient (only store best and current weights), simple, no hyper parameter (but needs validation data)
- Linear regression: input noise (leads to more robust solutions) is the same as L2

4 Nearest Neighbours, Trees and Forests

- Non-parametric methods: use training data directly for prediction (complexity adapts to training data, very fast training, slow predictions, hard for high dimensions)
- KNN: needs lots of training data and less than 20 attributes, can learn complex functions, regression works similar
- Increasing k reduces variance, increases bias
- Euclidean distance when each variable has same unit, otherwise normalize data
- Cosine Distance (documents, images), Hamming Distance (string data/categorical features), Manhattan Distance (coordinate-wise), Mahalanobis Distance (unaffected by coordinate transformations)
- Performance of KNN degrades with irrelevant dimensions/high dimensions (most points far away) \rightarrow dimensionality reduction, feature selection
- KD-Tree to find neighbours; Build: choose dimension by longest hyperrectangle side, median as pivot; Traverse: move down tree, find region containing \mathbf{x} , find closest \mathbf{x}^* , move up for regions intersecting hypersphere, update \mathbf{x}^*
- Regression/Classification Tree: split data into two at each node using criterion
- Splitting criterion regression: Minimum residual sum of squares: $RSS = \sum_{\text{left}} (y_i - \bar{y}_L)^2 + \sum_{\text{right}} (y_i - \bar{y}_R)^2$, \bar{y} : average label subtree (variance in subtrees minimized)
- Criterion classification: Minimum entropy: score = $N_L H(p_L) + N_R H(p_R)$, $H(p_i) = -\sum_k p_i(k) \log p_i(k)$: subtree entropy, $p_i(k)$: proportion of class k in subtree i
- Stop: Minimum number of samples per node / maximum depth (tree complexity)
- Trees: easy to compute, no distributional assumption, non-linear, automatic variable selection, easy to interpret, lower accuracy, sensitive to data change
- Random Forests: use multiple trees to improve accuracy
- Bagging: Fit trees to bootstrap samples from data (combine by voting/averaging)
- Ideal: linear variance reduction (trees correlated \rightarrow reduction still significant)
- Bagging: less variance, bias unaffected \rightarrow use strong trees (high variance/low bias)
- Random Forests: Also randomize considered variables at each splitting criterion, grow to maximum depth (loss of interpretability, good accuracy, less unstable)

5 Dimensionality Reduction and Clustering

- Motivation: Invert $\mathbf{X}^\top \mathbf{X}$ for linear regression: $d \times d \rightarrow O(d^3) \rightarrow$ find $d_{\text{new}} \ll d$
- Find (linear) mapping $\mathbf{x}_i \rightarrow \mathbf{z}_i$ to lower dimension with $\mathbf{z}_i = \mathbf{W} \mathbf{x}_i$
- Orthonormal basis system: $\mathbf{x} = \sum_i^D z_i \mathbf{u}_i \rightarrow z_i = \mathbf{u}_i^\top \mathbf{x} \rightarrow$ only use subset for dimensionality reduction (minimize squared reproduction error $\sum \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2$)
- Minimizing error \Leftrightarrow maximizing variance of projection (with zero mean data)
- Principle component analysis: find principal directions \mathbf{u}_i and their variance λ_i
- $\mathbf{u}_1 = \operatorname{argmax}_{\mathbf{u}} \frac{1}{N} \sum (\mathbf{u}^\top (\mathbf{x}_i - \boldsymbol{\mu}))^2$ s.t. $\mathbf{u}^\top \mathbf{u} = 1$, \mathbf{u}_2 maximizes variance in orthogonal complement of \mathbf{u}_1
- Objective can be written in terms of sample covariance: $E(\mathbf{u}) = \mathbf{u}^\top \boldsymbol{\Sigma} \mathbf{u}$
- Constraint Optimization: Lagrangian Multipliers ($L = \text{objective} - \text{multiplier} \cdot \text{constraint}$): $\min_x x^2$ s.t. $x \geq b \rightarrow \min_x \max_\alpha L(x, \alpha) = x^2 - \alpha(x - b)$ s.t. $\alpha \geq 0$ (Min forces max to behave such that constraints are satisfied)
- Dual formulation: $\boldsymbol{\lambda}^* = \operatorname{argmax} g(\boldsymbol{\lambda})$, $g(\boldsymbol{\lambda}) = \min_x L(x, \boldsymbol{\lambda})$ s.t. $\lambda_i \geq 0$, $\mathbf{x}^* = \operatorname{argmin}_x L(x, \boldsymbol{\lambda}^*)$ (swap min/max)
- Slaters condition: convex objective/constraints \Rightarrow dual \Leftrightarrow primal (original)
- PCA: $\mathbf{u}_1 = \operatorname{argmax}_{\mathbf{u}} \mathbf{u}^\top \boldsymbol{\Sigma} \mathbf{u}$ s.t. $\mathbf{u}^\top \mathbf{u} = 1 \Rightarrow L(\mathbf{u}, \lambda) = \mathbf{u}^\top \boldsymbol{\Sigma} \mathbf{u} + \lambda(\mathbf{u}^\top \mathbf{u} - 1) \Rightarrow \boldsymbol{\Sigma} \mathbf{u} = \lambda \mathbf{u}$ (eigenvalue problem, largest value: maximum variance, vector: direction)
- Representation has minimum MSE of all linear representations of same dimension
- PCA: Subtract mean, (normalize variance of each dimension), choose first M largest eigenvalues/their vectors of $\boldsymbol{\Sigma}$, $\mathbf{z}_i = \mathbf{B}^\top (\mathbf{x}_i - \boldsymbol{\mu})$, reprojection: $\tilde{\mathbf{x}}_i = \boldsymbol{\mu} + \mathbf{B} \mathbf{z}_i$
- Choose M : based on application performance/based on captured variance
- Applications: face detection, morphing, natural image patches, ...
- Clustering Group data using similarity measure ($D(A, B) = D(B, A)$, $D(A, B) = 0 \Leftrightarrow A = B$, $D(A, B) \leq D(A, C) + D(B, C)$)
- Hierarchical Clustering: Dendrogram (similarity: height of lowest shared node)
- Outlier: single isolated branch
- Heuristic search of all possible trees: Bottom-up / Top-down (find best division)
- Bottom-up: each sample in own cluster, merge closest two clusters, until single cluster left (requires distance measure for samples and clusters)

- Cluster similarity: single linkage (minimum distance between two points) / complete linkage (maximum distance) / average linkage / centroid linkage
- Hierarchical: any number of clusters, $O(n^2)$, local optima, subjective interpretation
- Flat Clustering: K-Means: minimize quantization error (sum of squared distances)
 $SSD(C, D) = \sum d(\mathbf{x}_i, c(\mathbf{x}_i))^2$
- Iteration: 1. pick K random centroids c_i , 2. assign each point to closest c_i , 3. move centroids to mean of assigned points, 4. go to step 2 until no change
- $SSD = \sum_i \sum_k \delta_{ik} d(\mathbf{x}_i, \mathbf{c}_k)^2 \rightarrow$ assignment minimizes w.r.t. δ_{ik} , adjustment w.r.t. \mathbf{c}_k
- K-Means locally minimizes SSD (depends on initialization, global NP-hard)
- K-Means++: first centroid random, each following centroid furthest from all others
- Choose K : objective function decrease on holdout set or Knee-finding method
- Knee-finding: plot SSD for K , pick point where decrease is no longer steep
- K-Means: converges quickly, local optima, not applicable to categorical data/noisy data/outliers, clusters must be convex

6 Density Estimation and Expectation Maximization

- Non-parametric models (don't know form of class-conditional density) \rightarrow estimate directly from data (histograms, kernel density, KNN)
- Histograms: general, need exponential data (curse of dimensionality), fixed region size, $p(\mathbf{x}) \approx \frac{K}{NV}$ (K points in region R , N : total points, V : volume of R)
- Center R on \mathbf{x} : Kernel density: fix V , determine K , KNN: fix K , determine V
- Kernel Density Estimation: $k(\mathbf{x}, \mathbf{y})$: non-negative, distance-dependent: $k(\mathbf{x}, \mathbf{y}) = g(\mathbf{x} - \mathbf{y})$, $V = \int g(\mathbf{u}) d\mathbf{u}$, $K(\mathbf{x}_*) = \sum g(\mathbf{x}_* - \mathbf{x}_i) \rightarrow p(\mathbf{x}_*) \approx \frac{K(\mathbf{x}_*)}{NV}$
- Parzen Window (hypercubes): $g(\mathbf{u}) = 1$ if $|u_j| \leq \frac{h}{2}, j = 1 \dots d$, else 0 $\Rightarrow p(\mathbf{x}_*) \approx \frac{K(\mathbf{x}_*)}{Nh^d}$, h : bandwidth, d : dimensionality (easy to compute, not very smooth)
- Gaussian Kernel: $g(\mathbf{u}) = \exp(-\frac{\|\mathbf{u}\|^2}{2h}) \rightarrow p(\mathbf{x}_*) \approx \frac{1}{N\sqrt{(2\pi h)^d}} \sum \exp(-\frac{\|\mathbf{x}_* - \mathbf{x}_i\|^2}{2h})$
(smooth, infinite support, computationally intensive, bigger $h \rightarrow$ smoother curve)
- Cross-validation for bin size/bandwidth/neighbours (highest likelihood on test-set)
- Mixture Models generality of non-parametric models and efficiency of parametric models \rightarrow create complex distribution by combining simple ones (e.g. Gaussians)
- Mixture coefficient \cdot component: $p(\mathbf{x}) = \sum p(k)p(\mathbf{x}|k) \sim$ any smooth density

- $p(k) = \pi_k \geq 0, \sum \pi_k = 1, p(\mathbf{x}|k) = N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \boldsymbol{\theta} = \{\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \pi_K, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K\}$
- Gradient descent on marginal log-likelihood? \rightarrow possible, inefficient (depends on all components, no closed form, slow convergence, sum does not go well with log)
- Mixture models \rightarrow latent variable models (observed variables \mathbf{x} and latent variables z): $p(\mathbf{x}, z|\boldsymbol{\theta})$ (parametric model), $p(\mathbf{x}|\boldsymbol{\theta}) = \sum_z p(\mathbf{x}, z|\boldsymbol{\theta})$ (marginal distribution)
- Kullback-Leibler Divergence (similarity of distributions): $KL(q||p) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$ (non-negative, zero for same distribution, non-symmetric \rightarrow no distance metric)
- Expectation-Maximization algorithm estimates latent variable models (iteratively increases lower bound of the marginal log-likelihood) \rightarrow local optima
- $\log p(\mathbf{x}|\boldsymbol{\theta}) = \sum_z q(z) \log \frac{p(\mathbf{x}, z|\boldsymbol{\theta})}{q(z)} + \sum_z q(z) \log \frac{q(z)}{p(z|\mathbf{x})} = L(q, \boldsymbol{\theta}) + KL(q(z)||p(z|\mathbf{x}))$ (decomposition holds for any $q(z)$, it makes optimization much simpler)
- $L(q, \boldsymbol{\theta}) \leq \log p(\mathbf{x}|\boldsymbol{\theta})$ contains joint distribution \rightarrow easier to optimize (often convex)
- Expectation step: Find $q(z)$ to minimize KL $\rightarrow q(z) = p(z|\mathbf{x}, \boldsymbol{\theta}_{\text{old}}) = \frac{p(\mathbf{x}, z|\boldsymbol{\theta}_{\text{old}})}{\sum_z p(\mathbf{x}, z|\boldsymbol{\theta}_{\text{old}})}$ (closed form for discrete z) \Rightarrow KL = 0 \Rightarrow lower bound $L(q, \boldsymbol{\theta}_{\text{old}})$ tight at $\boldsymbol{\theta}_{\text{old}}$
- Maximization step: Maximize $L(q, \boldsymbol{\theta})$: $\boldsymbol{\theta}_{\text{new}} = \text{argmax} \sum_z q(z) \log p(\mathbf{x}, z|\boldsymbol{\theta}) + \text{const}$
- Full dataset: $L(q, \boldsymbol{\theta}) = \sum_i (\sum_k q_{ik} \log p(\mathbf{x}_i, k|\boldsymbol{\theta}) - \sum_k q_{ik} \log q_{ik})$ with $q_{ik} = q_i(z = k)$ (one latent variable per data-point)
- Gaussian mixture models E-step: Compute “responsibilities” of components: $q_{ik} = \frac{\pi_k N(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j N(\mathbf{x}_i|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = p(z = k|\mathbf{x}_i)$
- M-step: Separate updates for additive objectives: $\boldsymbol{\pi} = \text{argmax} \sum_i \sum_k q_{ik} \log \pi_k$, $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k = \text{argmax} \sum_i q_{ik} \log N(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \rightarrow$ weighted ML estimation
- $\pi_k = \frac{\sum_i q_{ik}}{N}$, $\boldsymbol{\mu}_k = \frac{\sum_i q_{ik} \mathbf{x}_i}{\sum_i q_{ik}}$, $\boldsymbol{\Sigma}_k = \frac{\sum_i q_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top}{\sum_i q_{ik}}$
- EM for GMMs: Initialize (K-Means for component means and fixed covariance), until convergence: E-step (responsibilities q_{ik}), M-step (update $\boldsymbol{\pi}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$)
- EM very sensitive to initialization, K-Means special case of EM
- More components \rightarrow better likelihood (beware of overfitting)
- EM for dimensionality reduction (probabilistic PCA): $\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}$, (latent variable \mathbf{z} : low dimensional representation, $\boldsymbol{\mu}$: constant offset, $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I})$: noise)
- Continuous latent variable: $p(\mathbf{z}) = N(\mathbf{0}, \mathbf{I})$, observation model $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) = N(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ with parameters $\boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\mu}, \sigma^2\}$
- Generative process: sample $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$, project: $\mathbf{y} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu}$, add noise: $\mathbf{x} = \mathbf{y} + \boldsymbol{\varepsilon}$

- Maximize marginal loglike(θ) = $\sum_i \log(\int_z N(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) N(\mathbf{z}|\mathbf{0}, \mathbf{I}) d\mathbf{z}) \rightarrow \text{EM}$
- E-step: Posterior $q_i(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}_i, \theta)$ with $\boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}_i} = (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^\top (\mathbf{x}_i - \boldsymbol{\mu})$, $\boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}_i} = \sigma^2 (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I})^{-1}$ (only possible because \mathbf{x} is linear in \mathbf{z})
- M-step: $L(q, \theta) = \sum_i \mathbb{E}_{q_i(\mathbf{z})} [\log p(\mathbf{x}_i|\mathbf{z}, \theta)] + \text{const} \approx \sum_i \log p(\mathbf{x}_i|\mathbf{z}_i, \theta)$ with $\mathbf{z}_i \sim q_i(\mathbf{z})$ (approximate with single sample per \mathbf{x}_i) \rightarrow solution: standard least squares:

$$\begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{W} \end{bmatrix} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}, \sigma^2 = \frac{1}{nd} \sum_i \sum_k^d (y_{ik} - x_{ik})^2$$
- PCA with eigenvectors preferred (one step \rightarrow very fast), probabilistic PCA provides density, helps understand EM and more complex dimensionality reduction methods
- EM: assumes KL can be zero (posterior can be evaluated analytically), \mathbf{z} must be discrete/linear gaussian \rightarrow Variational Bayes/Inference can work with KL > 0

7 Kernel Methods

- Kernel: represent $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ by $[\mathbf{K}]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ ($k : X \times X \rightarrow \mathbb{R}$: comparison)
- Modularity between choice of k and algorithm, poor scalability
- Positive definite kernel function k : symmetric, \mathbf{K} is always positive definite
- $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$: positive definite kernel (arbitrary feature function ϕ)
- Theorem: positive definite (p.d.) kernel \Leftrightarrow associated feature space
- Kernel for polynomial features of degree d : $k(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle^d$
- Gaussian kernel: $k(\mathbf{x}, \mathbf{y}) = \exp(-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2\sigma^2})$ with bandwidth σ (most used kernel)
- Gaussian kernel is inner product of two infinite dimensional feature vectors \rightarrow p.d.
- Kernel trick: feature based algorithms can use infinite dimensional feature space if rewritten to contain inner products of feature vectors \rightarrow better than linear features
- Kernel ridge regression: $\mathbf{w}^* = (\Phi^\top \Phi + \lambda \mathbf{I})^{-1} \Phi^\top \mathbf{y} = \Phi^\top (\Phi \Phi^\top + \lambda \mathbf{I})^{-1} \mathbf{y} \rightarrow d \times d$ matrix inversion (infinite) to $N \times N$ matrix inversion (\mathbf{K} , by using matrix identity)
- \mathbf{w}^* still d -dimensional, but can evaluate $f(\mathbf{x}) = \phi(\mathbf{x})^\top \mathbf{w}^* = \mathbf{k}(\mathbf{x})^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y} = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x})$ with $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$
- Comparison to linear regression with gaussian features: kernel allows setting centers adaptively (fixed without kernel trick)
- Choose hyper-parameter bandwidth via cross-validation
- Kernel methods (ridge regression, gaussian processes, SVMs): must store all samples, high computation, flexible representation, good for small data, hard to scale

8 Support Vector Machines

- Classification: class labels 1 and -1 for SVMs ($f(\mathbf{x}_i)y_i > 0$)
- Scalar projection of \mathbf{a} on \mathbf{b} : $a_b = \frac{\mathbf{a}^\top \mathbf{b}}{\|\mathbf{b}\|}$
- Support vectors: data points closest to decision boundary (other samples ignored), maximize margin ρ
- Maximum margin classifier has smaller complexity \Rightarrow generalizes better
- Distance between point \mathbf{x}_i and line: $r = \frac{\mathbf{w}^\top \mathbf{x}_i + b}{\|\mathbf{w}\|}$
- Choose scaling for \mathbf{w}, b so that $\mathbf{w}^\top \mathbf{x}_+ + b = 1$ (positive support vector), $\mathbf{w}^\top \mathbf{x}_- + b = -1$ (negative support vector) $\rightarrow \rho = \frac{2}{\|\mathbf{w}\|}$
- Optimization problem: $\operatorname{argmax}_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$ s.t. $\mathbf{w}^\top \mathbf{x}_i + b \geq 1$ if $y_i = 1, \leq -1$ if $y_i = -1$ (one positive + negative point satisfy equality, otherwise weight could be reduced)
- Reformulation: $\operatorname{argmin}_{\mathbf{w}} \|\mathbf{w}\|^2$, s.t. $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1$ (convex, single optimum)
- Choose trade-off between margin and accuracy (for outliers) \rightarrow slack-variables $\xi_i \geq 0$ allow violation of margin: $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \Rightarrow \operatorname{argmin}_{\mathbf{w}, \xi} \|\mathbf{w}\|^2 + C \sum \xi_i$ (C : inverse regularization, small \rightarrow large ρ , large \rightarrow small ρ , infinite \rightarrow hard ρ)
- Reformulated as unconstrained optimization: $\operatorname{argmin}_{\mathbf{w}} \|\mathbf{w}\|^2 + C \sum_i \max(0, 1 - y_i f(\mathbf{x}_i)) = \text{regularization} + \text{hinge loss} \rightarrow$ convex, one minimum, but not differentiable \rightarrow similar to logistic regression loss
- Hinge loss: $\max\{0, 1 - y_i f(\mathbf{x}_i)\}$ / Logistic loss: $\log(1 + \exp(-y_i f(\mathbf{x}_i))) \rightarrow y_i f(\mathbf{x}_i)$ should be large for both / saturates if it gets too large
- Sub-gradient: Any \mathbf{g} at point \mathbf{x} so that $f(\mathbf{z}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{z} - \mathbf{x})$, if f is differentiable at $\mathbf{x} \Rightarrow \mathbf{g} = \nabla f(\mathbf{x})$
- Let $f(x) = \max\{f_1(x), f_2(x)\}$. $f_1(x) = f_2(x) \Rightarrow \mathbf{g} \in [\nabla f_1(x), \nabla f_2(x)]$
- Sub-gradient descent: $\mathbf{x}_{t+1} = \mathbf{x}_t + \eta \mathbf{g}$ (does not always decrease f , store best \mathbf{x}^*)
SVMs: each iteration, pick random (\mathbf{x}_i, y_i) . $y_i f(\mathbf{x}_i) < 1 : \mathbf{w}_{t+1} = \mathbf{w}_t - \eta(2\mathbf{w}_t - Cy_i \mathbf{x}_i)$, otherwise $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta 2\mathbf{w}_t$
- SVM: classification standard in 90s/00s (pedestrian detection, text categorization, character recognition, bioinformatics), extends to regression, outperformed by NNs
- Kernel SVM: Dual derivation: $\mathbf{w}^* = \sum_i \lambda_i y_i \phi(\mathbf{x}_i)$ (λ_i : constraint coefficient)
- $\frac{\partial L}{\partial b} = -\sum_i \lambda_i y_i \Rightarrow \sum_i \lambda_i y_i = 0$ no solution for b , but additional condition (b can be computed from \mathbf{w} : $b = y_i - \mathbf{w}^\top \phi(\mathbf{x}_i)$ (for \mathbf{x}_i on margin))

- Kernel trick for SVMs: $g(\boldsymbol{\lambda}) = \sum \lambda_i - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j)$
- Dual optimization (slack variables): $\max_{\boldsymbol{\lambda}} \sum \lambda_i - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j)$ s.t. $C \geq \lambda_i \geq 0$, $\sum \lambda_i y_i = 0$ with $b = y_k - \sum_i y_i \lambda_i k(\mathbf{x}_i, \mathbf{x}_k)$ where $C > \lambda_k > 0$ and $f(\mathbf{x}) = \sum_i y_i \lambda_i k(\mathbf{x}_i, \mathbf{x}) + b$ (upper bound C limits λ_i so misclassifications allowed)
- Control overfitting: set C (low $C \rightarrow$ low complexity), choose kernel, vary bandwidth

9 Bayesian Machine Learning

- Estimate $\boldsymbol{\theta}^*$ uncertainty, infinite predictors (mean) \rightarrow give prediction uncertainty
- Compute posterior $p(\boldsymbol{\theta}|D) = \frac{p(D|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(D)}$, $p(D|\boldsymbol{\theta})$: data likelihood, $p(\boldsymbol{\theta})$: prior (subjective belief), $p(D)$: evidence (normalization, later used for model comparison)
- Compute predictive distribution (marginal likelihood) $p(\mathbf{x}^*|D) = \int p(\mathbf{x}^*|\boldsymbol{\theta})p(\boldsymbol{\theta}|D)d\boldsymbol{\theta}$, $p(\mathbf{x}^*|\boldsymbol{\theta})$: likelihood (weighted ensemble method, often uses samples of $p(\boldsymbol{\theta}|D)$)
- Prior should express belief and domain knowledge $\xrightarrow{\text{ML}}$ weights should be small $\rightarrow p(\boldsymbol{\theta}) = N(\boldsymbol{\theta}|\mathbf{0}, \lambda^{-1}\mathbf{I})$, λ : precision of the prior
- Completing the square: Bring exponent in canonical squared form: $\exp(-\frac{1}{2}a\mu^2 + b\mu + \text{const}) \rightarrow$ for gaussian distributions: $\mu_N = a^{-1}b$, $\sigma_N^2 = a^{-1}$
- Posterior if prior/likelihood are gaussian: $\mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\mu_{\text{ML}} + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0$ and $\sigma_N^2 = \frac{\sigma^2\sigma_0^2}{N\sigma_0^2 + \sigma^2}$ with $\mu_{\text{ML}} = \frac{\sum x_i}{N}$ and μ_0, σ_0 from prior (variance decreases with more training samples, posterior interpolates between prior mean/sample average)
- Gaussian Propagation: predictive distribution is gaussian with $\mu_{x^*} = \mu_N$ and $\sigma_{x^*}^2 = \sigma_N^2 + \sigma^2$
- Conjugate prior for likelihood function \Leftrightarrow posterior/prior: same distribution family
- Bayesian Learning: For large datasets: point estimate, advantage for small dataset
- Simplification: Maximum a-posteriori solution: $\boldsymbol{\theta}_{\text{MAP}} = \text{argmax}_{\boldsymbol{\theta}} \log p(D|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$ maximizes posterior \rightarrow use for prediction: $p(\mathbf{x}^*|D) \approx p(\mathbf{x}^*|\boldsymbol{\theta}_{\text{MAP}})$
- MAP regression: gaussian prior \leftrightarrow L2; gaussian likelihood \leftrightarrow squared loss \leftrightarrow ridge regression ($\lambda_{\text{ridge}} = \lambda\sigma^2$, uncertainty only depends on estimated noise σ^2)
- Bayesian linear regression: likelihood: $p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = N(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I})$ (multivariate distribution, σ^2 : noise variance), prior: $p(\mathbf{w}) = N(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I})$
- Gaussian Bayes Rule 1/2 for evaluating $p(\mathbf{x}|\mathbf{y})$ (different derivations of posterior distribution): rule 1 if $\dim(\mathbf{y}) < \dim(\mathbf{x})$ (infinite dimensional features), otherwise rule 2

- Gaussian Propagation to evaluate $p(\mathbf{y})$
- Posterior/predictive mean equivalent to MAP estimate, but we get uncertainty for parameters: $\Sigma_{w|X,y} = \sigma_y^2(\Phi^\top \Phi + \sigma_y^2 \lambda \mathbf{I})^{-1}$ and variance $\sigma^2(\mathbf{x}^*)$ is input dependent
- Compute predictive distribution with gaussian propagation
- Gaussian Process: distribution over functions $f(\mathbf{x})$ so that any set \mathbf{t} of function values evaluated at $\mathbf{x}_1, \dots, \mathbf{x}_n$ is jointly gaussian distributed: $f(\mathbf{x}) \sim GP(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$, $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$: mean (prior belief about function, zero for simplicity), $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f(\mathbf{x})f(\mathbf{x}')]$: positive definite correlation of function evaluations at \mathbf{x}, \mathbf{x}'
- $p(\mathbf{t}|\mathbf{X}) = N(\mathbf{t}|\mathbf{0}, \mathbf{K}) \xrightarrow{\text{noise}} p(\mathbf{y}|\mathbf{X}) = N(\mathbf{y}|\mathbf{0}, \mathbf{K} + \sigma_y^2 \mathbf{I})$ ($y_i = f(\mathbf{x}_i) + \varepsilon$)
- Predictive: $\mu(\mathbf{x}^*) = \mathbf{k}(\mathbf{x}^*)^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}$, $\sigma(\mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{x}^*) + \sigma_y - \mathbf{k}(\mathbf{x}^*)^\top (\mathbf{K} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}^*)$ (mean \leftrightarrow kernel ridge regression + input dependent variance estimate (small for high kernel activations)) \rightarrow kernel-version of bayesian linear regression
- Weight space view: Bayesian: subsume prior precision λ into kernel (vector/matrix): $\mathbf{K} = \lambda^{-1} \Phi_X \Phi_X^\top$
- Function view (from Gaussian process) vs. weight space view (from Bayesian Linear Regression with kernel trick)
- Posterior derived from Bayesian view (rule 1): $\mu_{w|X,y}, \Sigma_{w|X,y}$ (potentially infinite dimensions) \rightarrow can evaluate predictive distribution with kernel trick (same as GP)
- GP: computationally hard ($O(N^3)$), very principled approach to regression learning
- Kernel parameters: weight precision λ , observation noise σ_y , length scale l (can be different per dimension): Gaussian Kernel: $k(\mathbf{x}_i, \mathbf{x}_j) = \lambda^{-1} \exp(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|}{2l^2}) + \delta_{ij} \sigma_y^2$
- Optimization: non-convex log-likelihood of data \rightarrow gradient descent (overfitting)
- GP: non-parametric Bayesian approach, prediction equations in closed form (gaussian), hyperparameter optimization complex, outperforms NNs for small datasets

10 Neural Networks

- Artificial neuron: $y = \phi(\mathbf{w}^\top \mathbf{x} + b)$ (like logistic regression)
- Feedforward network: directed acyclic graph (units grouped into layers)
- Fully connected layer (N inputs to M outputs): $\mathbf{y} = \phi(\mathbf{W}\mathbf{x} + \mathbf{b})$, $\mathbf{W} \in \mathbb{R}^{M \times N}$
- Activations: σ (0 to 1, kills gradient, not zero-centered (important for initialization), exp computationally hard), \tanh (-1 to 1, zero centered, kills gradient), ReLU (fast computation/convergence, not zero centered, $x < 0$: no gradient), leaky ReLU (fast), ELU ($\alpha(e^x - 1)$ for $x < 0$, benefits of ReLU, closer to zero mean, exp hard)

- Each layer computes function: $\mathbf{y} = f^L \circ \dots \circ f^1(\mathbf{x})$ (composite of functions)
- XOR: classic example why multiple layers are needed
- Linear layers \Leftrightarrow one linear layer \rightarrow need non-linearities \rightarrow FF-NNs can approximate any function (theoretically with single layer, but exponential number of units)
- Deterministic regression: $\mathbf{f} = \mathbf{W}^{(L)}\mathbf{h}^{(L-1)} + \mathbf{b}^{(L)}$, $l_i(\mathbf{x}_i, \boldsymbol{\theta}) = \text{squared loss}$ / Probabilistic: $p(\mathbf{y}|\mathbf{x}) = N(\mathbf{y}|\mathbf{W}^{(L)}\mathbf{h}^{(L-1)} + \mathbf{b}^{(L)}, \boldsymbol{\Sigma})$, $l_i(\mathbf{x}_i, \boldsymbol{\theta}) = -\log N(\mathbf{y}_i|\boldsymbol{\mu}(\mathbf{x}_i), \boldsymbol{\Sigma}(\mathbf{x}_i))$
- Deterministic classification: $\mathbf{f} = \mathbf{W}^{(L)}\mathbf{h}^{(L-1)} + \mathbf{b}^{(L)}$, $l_i(\mathbf{x}_i, \boldsymbol{\theta}) = \text{hinge loss}$ / Probabilistic: $\mathbf{f} = \sigma(\mathbf{W}^{(L)}\mathbf{h}^{(L-1)} + \mathbf{b}^{(L)})$, $l_i(\mathbf{x}_i, \boldsymbol{\theta}) = -c_i \log f(\mathbf{x}_i) - (1 - c_i) \log(1 - f(\mathbf{x}_i))$
- Deterministic multi-class: $\mathbf{f} = \mathbf{W}^{(L)}\mathbf{h}^{(L-1)} + \mathbf{b}^{(L)}$, loss not covered / Probabilistic: $\mathbf{f} = \text{softmax}(\mathbf{W}^{(L)}\mathbf{h}^{(L-1)} + \mathbf{b}^{(L)})$, $l_i(\mathbf{x}_i, \boldsymbol{\theta}) = -\sum_k \mathbf{h}_{c_i, k} \log f_k(\mathbf{x}_i)$
- NNs learn features that can be separated linearly by last layer
- Back-propagation learning algorithm: Compute $\frac{\partial L}{\partial \mathbf{W}^{(l)}}, \frac{\partial L}{\partial \mathbf{b}^{(l)}}$ recursively (chain rule)
- Computation graph: node: input, edge: node computed as function of other node
- Forward pass: compute loss / backward pass: compute derivatives
- Notation: $\bar{y} = \frac{\partial L}{\partial y}$ (error signals)
- Multivariate chain rule: $\frac{\partial}{\partial t} f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t} = \bar{x} \frac{\partial x}{\partial t} + \bar{y} \frac{\partial y}{\partial t} = \bar{t} \rightarrow$ in vector notation: $\frac{\partial}{\partial t} f(\mathbf{x}(t)) = \frac{\partial f}{\partial \mathbf{x}}^\top \frac{\partial \mathbf{x}}{\partial t}$
- Backpropagation (v_1, \dots, v_N in topological order): $\forall i$: Compute v_i as a function of $\text{Pa}(v_i)$ (forward pass) $\rightarrow \bar{v}_N = 1 \rightarrow \forall i : \bar{v}_i = \sum_{j \in \text{Ch}(v_i)} \bar{v}_j \frac{\partial v_j}{\partial v_i}$ (backward pass)
- Chain rule for matrix-vector products: $\nabla_{\mathbf{W}} f = \frac{\partial f(\mathbf{z})}{\partial \mathbf{W}} = \frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} \frac{\partial}{\partial \mathbf{W}} (\mathbf{W}\mathbf{x} + \mathbf{b}) = \frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} \mathbf{x}^\top$ ($\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$)
- Forward: one add-multiply operation per weight / backward: two times forward (cost linear in number layers, quadratic in units per layer)
- Backprop neurally implausible (biological alternatives much slower on computers)
- Problems with standard SGD: slow along shallow dimension, jitter along steep dimension; stuck in local minima; noisy loss function due to mini-batches
- Momentum term (running gradient average): $\mathbf{m}_{k+1} = \gamma_k \mathbf{m}_k + (1 - \gamma_k) \nabla L$ / Geometric Average: $\mathbf{m}_k = (1 - \gamma) \sum_i \gamma^{k-i} \mathbf{g}_i$
- Gradient normalization (RMSProp): large steps in plateaus, small steps in steep areas: $\mathbf{g}_k = \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_k)$, $\mathbf{b}_{k+1, i} = \gamma \mathbf{v}_{k, i} + (1 - \gamma) \mathbf{g}_{k, i}^2$, $\boldsymbol{\theta}_{k+1, i} = \boldsymbol{\theta}_{k, i} - \frac{\eta}{\sqrt{\mathbf{b}_{k+1, i} + \varepsilon}} \mathbf{g}_{k, i}$ per dimension i , \mathbf{v}_k : average of gradient norms, ε : prevent division by zero

- Adam: adaptive momentum + normalization: $\theta_{k+1,i} = \theta_{k,i} - \frac{\eta c_2(k)}{\sqrt{c_1(k)v_{k+1,i} + \varepsilon}} m_{k+1,i}$
(no convergence guarantee, underestimation at start fixed by $c_i(k) = \frac{1}{1-\gamma_i^k}$)
- Learning rate decay: Reduce at fixed points / Cosine: $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(\frac{t\pi}{T}))$ /
Linear $\alpha_t = \alpha_0(1 - \frac{t}{T})$ / Inverse root: $\alpha_t = \frac{\alpha_0}{\sqrt{t}}$ (T : total number of epochs)
- First order optimization: step in direction of minimum of linear approximation /
second order optimization: step to minimum of quadratic approximation
- $\theta^* = \theta_0 - \frac{1}{2}\mathbf{H}^{-1}\mathbf{g}$ with Hessian $H = \nabla_{\theta}^2 L(\theta)$ (no hyperparameters, no learning
rate, less iterations, inverse in $O(N^3)$, N is in the millions)
- Solutions: quasi-Newton methods (BFGS, approximate Hessian over time) / Limi-
ted memory BFGS (does not store full \mathbf{H}^{-1} , works well in full batch) \rightarrow in practice
use Adam or L-BFGS (only on full batch with small noise)
- Regularization: model ensembles: train multiple models (or use snapshots of one
during training), average their results
- Dropout: randomly (often 50%) set neurons to zero (in each forward pass) \rightarrow forces
redundancy, can be interpreted as ensembles with shared parameters
- Testing dropout: average over multiple dropout masks (ensemble view) / multiply
each weight by dropout rate (expectation view)
- Drop connect: drop neuron connections (training) / use all connections (testing)
- Data preprocessing: initialization optimized for zero-mean unit variance data /
PCA / whitening of low-d data (covariance matrix is \mathbf{I})
- Classification loss less sensitive to small weight changes after normalization
- Weight initialization: constant \rightarrow all gradients equal, no distinct features can be
learned \rightarrow random initialization needed
- Fixed variance \rightarrow activations go to zero/saturate over deep layers (no gradients)
- Xavier initialization ($\sigma_W = \frac{1}{\sqrt{D_{\text{in}}}}$): activations nicely scaled for all layers (for tanh)
/ For ReLU: $\sigma_W = \frac{2}{\sqrt{D_{\text{in}}}}$
- Practice tips
 1. Check initial loss (without L2 should be $\log C$ for softmax with C classes)
 2. Overfit small sample (get 100% training accuracy, change architecture / η)
 3. Find η to strongly decrease loss in 100 iterations (full training data, small L2)
 4. Grid search around η / L2 from previous step (train each for 1 to 5 epochs)
 5. Train best models from step 4 for longer (10 to 20 epochs) without η decay

- 6. Loss curves: plateau end $\rightarrow \eta$ decay / plateau beginning \rightarrow bad initialization / plateau after η step decay \rightarrow decay later / validation accuracy going up \rightarrow train longer / validation accuracy going down \rightarrow regularize/more data / same training/validation accuracy \rightarrow train longer/bigger model
- NNs work very well, even though we have more parameters than training samples

11 Convolutional and Recurrent Neural Networks

- Image inputs \rightarrow huge amount of weights with FC-layers
- Close pixels more correlated \rightarrow use convolutions (slide filter over image)
- Stack filters to obtain multi-channel output
- Stride S : step-size ($> 1 \rightarrow$ down-sampling) / (zero)-padding P : fill image borders
- Convolutional Layer: $W_1 \times H_1 \times D_1 \rightarrow W_2 \times H_2 \times D_2$ with $W_2 = (W_1 - F + 2P)/S + 1$, $H_2 = (H_1 - F + 2P)/S + 1$ and $D_2 = K$: number of filters, F : kernel size
- (Max)-Pooling: smaller output dimension (applied to each channel with $P = 0$)
- Convolutional network: Convolutional layers, activations, pooling, FC layers at end
- Optimize deep models: residual block computes $F(x) + x \rightarrow$ new layers do no harm with $F(x) = 0$ at beginning
- Transfer learning (for small datasets): Convolutional layers are generic \rightarrow reusable (only train last FC layer(s) and freeze rest)
- AlexNet (2012): first use of ReLU, 8 layers (first CNN winner of ImageNet) / VGG (2014): more smaller filters (more non-linearities, fewer parameters, 19 layers) / ResNet (2015): very deep using residual connections (152 layers)
- Recurrent NNs: one to many (image captioning) / many to one (sentiment) / many to many (translation, video classification) \rightarrow use old state as input
- State $\mathbf{h}_t = f_W(\mathbf{h}_{t-1}, \mathbf{x}_t) = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t)$, $\mathbf{y}_t = \mathbf{W}_{hy}\mathbf{h}_t$
- Computational graph: unroll time steps \rightarrow network depth T , reuse \mathbf{W} each step
- Backpropagation through time (BPTT): forward/backward through entire sequence
- Truncated BPTT: keep \mathbf{h} , but only backpropagate for smaller number of steps
- Image Captioning: $\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{ih}\mathbf{v})$ with \mathbf{x} : previous word, \mathbf{v} : last layer CNN
- $\nabla \mathbf{h}$ depends on \mathbf{W} : largest singular value > 1 : exploding gradients \rightarrow scale gradient / largest singular value < 1 : vanishing gradients \rightarrow different RNN architecture

- Long-term short-term memory (LSTM): gated contribution of state/input (forget (erase cell), input (write), g (how much to write), output (how much to reveal))
- LSTM: Backpropagation from c_t to c_{t-1} : only elementwise multiplication (uninterrupted gradient flow similar to ResNet)
- Stack: layer 1 output sequence \rightarrow layer 2 input (only dropout non-recurrent edges)
- Gated Recurrent Units: no explicit f gate, less parameters, similar performance

12 Wrap-Up

Chapter	Classical Supervised Learning	Classical Unsupervised Learning	Kernel Methods	Bayesian Learning	Neural Networks
Algorithms	<i>Regression:</i> Linear, Ridge, KNN, Trees, Forests <i>Classification:</i> Logistic Regression, KNN, Trees, Forests	(p)PCA Clustering: Agglomerative, K-Means, EM for GMMs Density Estimation: KDE, KNN, Mixture Models	Kernel Regression, SVMs	Bayesian Linear Regression, Gaussian Processes	FF-NNs, Backprop, CNNs, LSTMs
Basics	Matrix and Vector Calculus, Probability Theory, MLE, Gradient Descent	Constraint Optimization, EM	Sub-gradients, Constraint Optimization	“Completing the square”, Gaussian Conditioning	Most of the others
Representations	Features, Basis Functions, Instances, Trees	Instances, Linear Projections, Centroids, Mixture Models	Kernels	Features, Kernels	NNs
Optimization	Least-squares, Gradient Descent	Eigen-Value Decomposition, EM	Sub-gradients, Quadratic Solver	Computing the Posterior	Adam, 2nd Order Gradient, Sub-gradients (ReLU)
Loss	MSE/SSE, Gaussian Log-Likelihood, Binary Cross Entropy, Soft-Max Likelihood	Reproduction Error, SSD, Sum of Discrepancies, Marginal Log-Likelihood	Maximum Margin, Hinge Loss	MAP, Posterior Approximation	Most of the others