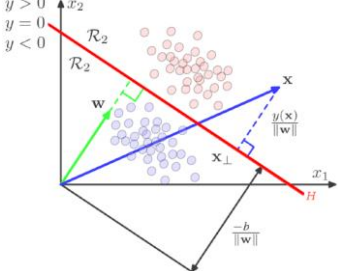
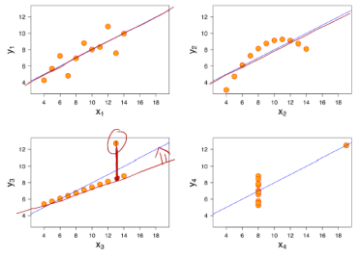
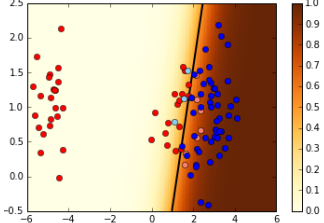
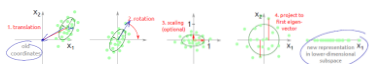
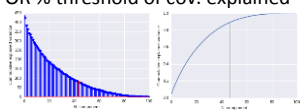
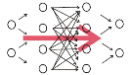
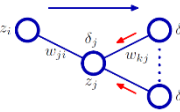
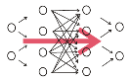

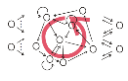
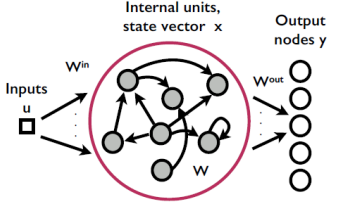


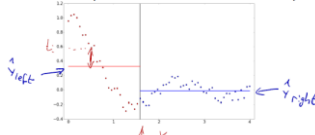
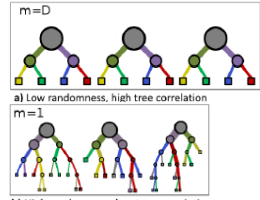
Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Linear Discriminant Analysis (LDA)</b>	ML2.x Bishop 4.x	<p><b>Decision function</b> <math>y(x) = w^T x + b</math> (generalized LDA with non-linear basis functions: <math>y(x) = w_0 + \sum_{i=1}^D w_i \phi(x)</math>)</p>  <p><b>Fishers Criterion:</b> maximize distance of projected means, minimize within-class variance <math>s_k^2</math> of projected classes k</p> $J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \xrightarrow{\text{explicit proj.}} \frac{w^T S_B w}{w^T S_W w}$ <p>for <math>\arg \max_w</math>: large scatter <math>S_B</math> “between-class” and low scatter <math>S_W</math> “within-class”</p> <p><b>analytic sol.:</b> <math>w = S_W^{-1}(m_2 - m_1)</math> and <math>b = -\frac{1}{2} w(m_1 + m_2)</math> with <math>S_W = S_{W1} + S_{W2}</math></p>	<p>weight vector <math>w \in R^{D \times 1}</math> (perpendicular to decision boundary <math>H</math>, projects data to one dim)</p> <p>bias/intercept <math>b</math> (or <math>w_0</math>) fixed offset</p> <hr/> <p><b>Free parameters:</b> <math>w: \frac{D(D+1)}{2}</math> (<math>w \in R^{D \times D}</math> is symmetric -&gt; calculate only half + diagonal; with intercept <math>w_0</math> it is <math>D + 1</math>)</p> <hr/> <p>A1: <b>Gaussian</b> data distributions for both classes (-&gt; each class can be described by covariance matrix)</p> <p>A2: <b>equal covariance matrices</b> of the classes</p>	<p>Supervised Learning</p> <p>Classification</p> <p>Train model with <b>labelled data</b> <math>\{(x_i, y_i)\}_{i=1}^N</math> to predict <b>discrete</b> label <math>y_{N+1}</math> from data <math>x_{N+1}</math></p> <hr/> <ul style="list-style-type: none"> <li>As a first shot method</li> <li>when noise model of sensors is known</li> <li>when class means are informative</li> </ul>	<p><b>training:</b> <math>O(D^3)</math> for matrix inv.  Strassen algorithm <math>O(D^{2.8})</math>  for extremely large <math>D</math> / not used in practice: Coppersmith-Winograd method <math>O(D^{2.3})</math></p> <p><b>test/recall:</b> <math>O(D)</math></p>	<p>+ weight vector <math>w</math> and bias <math>b</math> can be computed analytically (key: within-class covariance matrices are identical and Gaussian)</p> <p>- covariance matrix needs to be estimated and inverted</p>	<p>Various LDA <b>formulations</b> (gradient descent, eigenvalue problem, incremental for online learning)</p> <p>Generalize to <b>multiple classes</b>: one-against-rest (bad) one-against-one (better) inherent multiclass formulation: <math>y_k(x) = w_k^T x + b_k</math></p> <p><b>Easier condition:</b> <math>S_k = \sigma^2 I</math> (sphere) decision boundary perpendicular to <math>m_2 - m_1</math> search for nearest class mean</p> <p><b>Harder condition:</b> <math>S_k</math> arbitrary normal distr., different per class decision hyperquadratic (instead of hyperplane)</p>
<b>Linear Regression</b>	ML3.x Bishop 3.x	<p><b>Regression model</b> <math>y(x, w) = w^T x + \varepsilon</math> <math>y(x, w) = w_0 + w_1 x_1 + \dots + w_D x_D + \varepsilon</math> data point <math>x = (x_1, \dots, x_D)^T</math>, prediction <math>\hat{y} = w^T x</math>, residual error <math>\varepsilon</math></p> <p>(generalized with non-linear basis functions: <math>y(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \phi(x) + \varepsilon_j</math>)</p>  <p><b>Least Squares Loss Function:</b> Prediction <math>\hat{y}(x, w) = w^T X = Xw</math> <math>\arg \min_w \ \varepsilon\ ^2 = \arg \min_w \ y - \hat{y}\ ^2</math> <b>analytic sol.:</b> <math>X^T y = X^T X w</math> (“normal eq.”) <math>\Leftrightarrow w = (X^T X)^{-1} X^T y</math></p> <p><b>Regularization</b> combine Loss functions with penalty for large weight vectors <math>w</math></p>	<p>weight vector <math>w \in R^{D \times 1}</math> (single weight <math>w_i</math>: determines how sensitive prediction <math>\hat{y}</math> is to change of corresponding input var <math>x_i</math> [partial deriv. w.r.t <math>x_i</math>])</p> <p>bias/intercept <math>b</math> (or <math>w_0</math>) fixed offset with aug. notation, dim of data is <math>D + 1</math> with <math>x_0 = 1</math></p> <hr/> <p><b>Free parameters:</b> <math>w: M &gt; \frac{D(D+1)}{2}</math></p> <hr/> <p>A1: residuals means: <math>\forall i: E(\varepsilon_i) = 0</math> ≠ offset/skew/outliers/compare A1</p> <p>A2: Residuals uncorrelated and share same var: <math>\forall i: \text{Var}(\varepsilon_i) = \sigma^2</math> ≠ heteroscedastic data (e.g. % err.)</p> <p>A3: Residuals normally distributed <math>\varepsilon_i \sim N(0, \sigma^2)</math> (with zero mean -&gt; A1) ≠ offset/skew/outliers/compare A1</p>	<p>Supervised Learning</p> <p>Regression</p> <p>Train model with <b>labelled data</b> <math>\{(x_i, y_i)\}_{i=1}^N</math> to predict <b>continuous</b> label <math>y_{N+1}</math> from data <math>x_{N+1}</math></p> <hr/> <ul style="list-style-type: none"> <li>predict unknown values <math>y_i</math> for new <math>x_i</math></li> <li>estimate influence of single/several input var. (est. strength of corr. betw. <math>x_i</math> and <math>y</math>)</li> <li>Visualize relationships</li> </ul>	<p><b>training:</b> matrix inv <math>O(D^3)</math> (or <math>O(D^{2.3})</math>)</p> <p><b>test/recall:</b> <math>O(D)</math></p>	<p>+ useful analytical / computational properties</p> <p>- practical applicability limited by curse of dim / linearity</p>	


Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Logistic Regression</b>	ML4.x Bishop4.3.2	<p>wrap output of linear model <math>z = w^T x</math> with sigmoid/logistic function <math>g(z) = \frac{1}{1+e^{-z}}</math>:</p> <p><b>Model</b> <math>h_w(x) = g(w^T x) = \frac{1}{1+e^{-w^T x}}</math></p> <p>output is bound to <math>0 \leq h_w(x) \leq 1</math>, can be used as probability whether <math>x</math> belongs to one (<math>y = 1</math>) or the other class (<math>y = 0</math> or <math>-1</math>)</p>  <p><b>Logistic Cost Function</b>  <math>J(x) = \sum_i^N \log(\exp(-y_i(w^T x_i)) + 1)</math>  use log (quadratic loss <math>J(x) = \ h_w(x), y\ ^2</math> is non-convex because of sigmoid)</p> <p><i>name: formally derived as regression, but used as classification</i></p>	<p>weight vector <math>w \in R^{D \times 1}</math></p> <p>bias/intercept <math>b</math> (or <math>w_0</math>) fixed offset</p> <hr/> <p><b>Free parameters:</b> <math>w: M &gt; \frac{D(D+1)}{2}</math></p> <hr/> <p>A: same as LDA (?)</p>	<p>Supervised Learning</p> <p>Classification (!)</p> <hr/> <p>Train model with <b>labelled data</b> <math>\{(x_i, y_i)\}_{i=1}^N</math> to predict <b>discrete</b> label <math>y_{N+1}</math> from data <math>x_{N+1}</math></p> <hr/> <ul style="list-style-type: none"> <li>if probabilities of label are needed</li> </ul>	<p><b>training:</b> <math>O(D^3)</math></p> <p><b>test/recall:</b> <math>O(D)</math></p>	<p>+ Convex cost function</p> <p>– There is no analytic solution</p>	<p><b>Linear regression</b> for <math>g(x) = x</math> (identity)</p>
<b>Principal Component Analysis (PCA)</b> <p>singular value decomposition SVD eigenvalue decomposition EVD Karhunen-Loève transform</p>	ML5.x	<p>Project D-dim (<math>R^D</math>) data into a M-dim subspace (<math>R^M \subset R^D</math>) with <math>M \ll D</math>, that (1) contains the <b>relevant part</b> of our data and (2) such that data is <b>uncorrelated</b> (orthogonal basis vectors)</p>  <p><b>PCA-Transformation</b>  maximize variance <math>\sum_{i=1}^N (u_m^T x_i - u_m^T \bar{x})^2 = u_m^T S u_m + \lambda_m (1 - u_m^T u_m)</math>  (<math>S</math>: data covariance matrix, <math>\lambda_m</math>: Lagrange multiplier for unit vec constraint <math>u_m^T u_m = 1</math> to avoid <math>u_m \rightarrow \infty</math> for <math>\arg\max_{u_m} u_m^T S u_m</math>)  <b>analytic sol.:</b> deriv. w.r.t. <math>u_m</math> set to zero:  <math>S u_m = \lambda_m u_m</math> (Eigenvalue problem)  Take <math>M</math> eigenvectors <math>u_1, \dots, u_m</math> corresp. to the largest eigenvalues <math>\lambda_1, \dots, \lambda_m</math> (sort asc.)  e.g. until 99% of variance is preserved</p>	<p>eigenvalue spectrum <math>\langle (\lambda_m, u_m) \rangle_{m=1}^M</math> of data covariance matrix <math>S \in R^{D \times D}</math></p> <p><b>hyperparameter:</b>  <math>M</math>: nr. of output dim  OR % threshold of cov. explained</p>  <hr/> <p><b>Free parameters:</b> ...</p> <hr/> <p>A1: relevance is expressed by variance  A2: PCA determines a <b>linear</b> subspace</p>	<p>Unsupervised Learning</p> <p>linear feature extractor</p> <hr/> <p>pre-process:</p> <ul style="list-style-type: none"> <li>data compression</li> <li>dimensionality reduction</li> </ul>	<p><b>training:</b> <math>O(D^3)</math></p> <p>or iterative methods that scale better for high dim.</p> <p><b>test/recall:</b> <math>O(D)</math></p>	<p>+ popular</p> <p>– prone to outliers / scaling of a variable (cm instead of m)</p> <p>– doesn't guarantee good class-seperability (unsupervised; doesn't take class labels into acc)</p> <p>– computational complexity of calculating eigenvalues grows cubically in D</p>	<p>For <b>large dimensionality D</b>, using covariance matrix has disadvantages (<math>O(D^3)</math>).</p> <p>Various PCA <b>formulations</b> and (SVD, Bayesian PCA, iterative vs. analytical, Raleigh coeff., ...)</p> <p><b>Related subspace methods:</b>  Whitening / Sphering: transform data to zero mean and unit covariance (common preprocessing step)  Factor analysis (FA): incorporate domain-specific assumptions  Canonical correl. analysis (CCA): relate two data sources to a common subspace which maximizes cross-covariance  Kernel-PCA  non-linear extension of PCA</p>

Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Independent Component Analysis (ICA)</b>	MLSb.x	<p><b>Mixture/forward model</b> <math>x = As</math> with both unknown mixing matrix <math>A \in R^{N \times N}</math> and sources <math>s_1, \dots, s_n</math> maximize statistical independence (vs. PCA: max. cov.)</p> <p><b>backward model</b> <math>s = A^{-1}x = Wx</math> generative (describes how data <math>x</math> is generated by mixing process)</p> <p><b>Gradient descent (GD) on w:</b>  <math>\hat{s} = w^T x</math>  1: init <math>w</math>  2: determine direction in which kurtosis of <math>\hat{s}</math> decr./incr. most strongly (for pos/neg kurtosis)  3: run a step with GD for improved <math>w</math></p>	<p>...</p> <hr/> <p><b>Free parameters:</b> ...</p> <hr/> <p>A1: source signals <math>s_i</math> are independent of each other  A2: observed signals <math>x_i</math> are a linear mixing of source signals <math>s_i</math> with a fixed mixing matrix <math>A</math>.  A3: sources have a non-Gaussian distribution (one may be Gaussian).</p> <p>Helpful: mean-free data  Helpful: whitened data (<math>A \rightarrow \tilde{A}</math>)  reduce nr of parameters to be estimated from <math>D^2</math> to <math>\frac{D(D-1)}{2}</math> since</p>	<p>Unsupervised Learning</p> <hr/> <p>computational method for separating a multivariate signal (e.g. sounds) into (independent) additive subcomponents (e.g. speakers or singers)</p> <p>examples:  • audio sources (no reverberation, movement of sources)  • EEG sources</p>	<p><b>training:</b>  <math>O(D^3)</math></p>	<p>+ many variants (with varying interpretation of “statistical independence” for specific req of data)</p> <p>linear method:  + once trained can be applied extremely fast (e.g. online systems)  + independent components can be visualized</p> <p>– sources have arbitrary sign, order and amplitude (finding matching components is not trivial)</p>	<p>kurtosis hard to estimate robustly, other measures of non-Gaussianity: (max) negentropy, (min) mutual information, (max) likelihood, ...</p> <p><b>Related subspace methods:</b>  &gt; compare PCA  Blind Source Separation (BSS)  ICA is a special case of BSS</p>

Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Multilayer Perceptron (MLP)</b>  <b>Feed-Forward Neural Network (NN)</b>  	ML7.x Bishop 5	<p>Make basis functions linear comb. of input with adaptive parameters and train them:  Extension of logistic regression model  <math>y = \sigma(w^T \phi(x))</math>  <math>\sigma</math>: output activation func.  <math>\phi</math>: basis function with learned params <math>w^{(l)}</math>  Network organized in layers, output units of <math>l^{\text{th}}</math> layer serve as input of the <math>l + 1^{\text{th}}</math> layer</p> <p><b>network model:</b> <math>a_{ul} = \sum_{u_{l-1}=1}^{D_{l-1}} w_{ul u_{l-1}}^{(l)} z_{u_{l-1}}</math>  <math>z_{u_l} = h(a_{u_l})</math>: <math>u_l^{\text{th}}</math> unit of layer <math>l</math>  <math>D_l = \dim(z_{u_l})</math>: nr of hidden units of layer <math>l</math>  <math>u_l</math>: idx of units of layer <math>l</math>  <math>w_{u_l u_{l-1}}^{(l)}</math>: <math>D_l \times D_{l-1}</math> weights of layer <math>l</math></p> <p><math>z_{u_0} = x_i</math>: <math>i^{\text{th}}</math> input "unit"/data point  <math>D_0 = \dim(x_i)</math>: nr of input var. <math>x_i = (x_{i1}, \dots, x_{iD_0})</math>  <math>a_{u_2} = y_k</math>: <math>D_k</math>-Dimensional output in example three-layer network:  <math>y_k = \sigma \left( \sum_{k=0}^{D_2} w_{ku_2}^{(3)} h \left( \sum_{u_1=0}^{D_1} w_{u_2 u_1}^{(2)} h \left( \sum_{i=0}^{D_0} w_{u_1 i}^{(1)} x_i \right) \right) \right)</math></p> <p><b>generic MLP learning algorithm</b>  1: choose initial weight vector <math>w</math>  2: initialize minimization approach  3: <b>while</b> error did not converge <b>do</b>  4:   <b>for</b> all <math>(x_n, t_n) \in \mathcal{D}</math> <b>do</b>  5:     <math>y_n \leftarrow</math> forward pass(<math>x_n</math>)  6:     <math>\frac{\partial E_n}{\partial w_{u_l u_{l-1}}} \leftarrow</math> backward pass(<math>y_n, t_n</math>)  7:   <b>end for</b>  8:   sum over all n errors <math>\sum_n \frac{\partial E_n}{\partial w_{u_l u_{l-1}}}</math>  9:   perform update step  10: %one epoch: all data considered once%  11: <b>end while</b></p> <p><b>training:</b> with data <math>\mathcal{D} = \langle (x_n, t_n) \rangle_{n=1}^N</math>  <math>\text{argmin}_w [E(w, \mathcal{D}) = \sum_{n=1}^N E_n(w)]</math>  <b>backpropagation</b> (chain rule)  partial derivative of error <math>E_n</math> w.r.t weight <math>w_{u_l u_{l-1}}</math>:  <math>\frac{\partial E_n}{\partial w_{u_l u_{l-1}}} = \frac{\partial E_n}{\partial a_{u_l}} \frac{\partial a_{u_l}}{\partial w_{u_l u_{l-1}}} = \delta_{u_l} z_{u_{l-1}}</math></p> <p>for output unit <math>y_k</math>: <math>\delta_k = \frac{\partial E_n}{\partial a_k} = \frac{\partial E_n}{\partial y_k} \frac{\partial y_k}{\partial a_k} = y_k - t_k</math>  for hidden unit <math>z_{u_l}</math>: <math>\delta_{u_l} = \frac{\partial E_n}{\partial a_{u_l}} = \sum_{u_{l+1}} \frac{\partial E_n}{\partial a_{u_{l+1}}} \frac{\partial a_{u_{l+1}}}{\partial a_{u_l}} = \sum_{u_{l+1}} \frac{\partial E_n}{\partial a_{u_{l+1}}} \frac{\partial a_{u_{l+1}}}{\partial z_{u_l}} \frac{\partial z_{u_l}}{\partial a_{u_l}} = \sum_{u_{l+1}} \delta_{u_{l+1}} w_{u_{l+1} u_l}^{(l+1)} h'(a_{u_l})</math></p> 	<p><b>nonlin. Activation Functions <math>h(a_j)</math>:</b>  transform activations of lin. comb.  <math>h_{\text{linear}}(a_j) = a_j</math> <math>h'(a_j) = 1</math>  <math>h_{\text{logistic}}(a_j) = \frac{1}{1 + \exp(-a_j)}</math>  <math>h'(a_j) = h(a_j)(1 - h(a_j))</math>  <math>h_{\text{tanh}}(a_j) = \frac{e^{a_j} - e^{-a_j}}{e^{a_j} + e^{-a_j}}</math>  <math>h'(a_j) = 1 - h(a_j)^2</math>  <math>h_{\text{relu}}(a_j) = \max(0, a_j)</math>  <math>h'_{a_j &lt; 0}(a_j) = 0</math> <math>\forall h'_{a_j \geq 0}(a_j) = 1</math></p> <p><b>pairs of output act. and error func.:</b>  binary classification: <math>\sigma_{\text{logistic}}(a)</math>  prob. interpret. <math>0 \leq y(x, w) \leq 1</math>: <math>p(C_1 x)</math>  Cross-entropy error function (neg. log likelih.)  <math>E(w) = -\sum_n (t_n \ln y_n + (1 - t_n) \ln(1 - y_n))</math>  regression: <math>\sigma_{\text{linear}}(a)</math>  <math>E(w) = -\sum_{n=1}^N (y_n - t_n)^2</math>  multiclass classification:  generalization of cross-entropy err.:  <math>E(w) = -\sum_{n=1}^N \sum_{k=1}^K (t_{kn} \ln y_k(x_n, w))</math></p> <p><b>forward pass</b>  1: <b>for</b> all input units <math>i</math> <b>do</b>  2:   set <math>z_{u_0} = x_i(n)</math>  3: <b>end for</b>  4: <b>for</b> all hidden layers <math>l</math> (forward) <b>do</b>  5:   <b>for</b> all units <math>u_l</math> in layer <math>l</math> <b>do</b>  6:     set <math>a_{u_l} = \sum_{u_{l-1}=1}^{D_{l-1}} w_{u_l u_{l-1}}^{(l)} z_{u_{l-1}}</math>  7:     set <math>z_{u_l} = h(a_{u_l})</math>  8:   <b>end for</b>  9: <b>end for</b>  10: <b>for</b> all output units <math>k</math> <b>do</b>  11:   set <math>a_{u_L} = \sum_{k=0}^{D_k} w_{ku_L}^{(L)} z_{u_L}</math>  12:   set <math>y_k = \sigma(a_{u_L})</math>  13: <b>end for</b></p> <p><b>Free parameters:</b> <math>\sum_{l=1}^L D_l \times D_{l-1}</math>  <math>D_0 = D_{\text{in}} = \dim(x_i)</math> and <math>D_L = D_{\text{out}} = \dim(y_k)</math>  for two layers: <math>D_{\text{in}} \times D_1 + D_1 \times D_{\text{out}}</math></p> <p>A1: error function is continuous  A2: error function is differentiable</p>	<p>Supervised Learning</p> <p>Nonlinear Regression / Classification / Multiclass Classification  (<math>K</math> output units and softmax activation to <math>y_k = \frac{\exp(a_k(x, w))}{\sum_j \exp(a_j(x, w))}</math>)</p> <p>• Representation Learning:  Parametrize basis functions and adapt the data to discover useful representations automatically</p> <p><b>backward pass</b>  1: <b>for</b> all output units <math>k</math> <b>do</b>  2:   compute <math>\delta_k = y_k - t_k</math>  3: <b>end for</b>  4: <b>for</b> all hidden layers <math>l</math> (backward) <b>do</b>  5:   <b>for</b> all units <math>u_l</math> in layer <math>l</math> <b>do</b>  6:     backpropagate existing <math>\delta</math>'s to obtain <math>\delta_{u_l}</math>  7:     compute <math>\frac{\partial E_n}{\partial w_{u_l u_{l-1}}} = \delta_{u_l} z_{u_{l-1}}</math>  8:   <b>end for</b>  9: <b>end for</b></p>	<p><b>training:</b>  [feedforward &amp; backward pass] * nr. of gradient descent steps</p> <p><b>test/recall:</b>  feedforward pass</p>	<p>+ many applications resulting model can be significantly more compact (than SVM with same generalization prop) -&gt; faster to evaluate</p> <p>- price for compactness: no convex likelihood (optimization) function</p> <p>- not robust to hyperparameter settings</p>	<p><b>Perceptron</b>  uses step-function as non-linearities (hence nondifferentiable)</p> <hr/> <p><b>two layer example</b>  first (hidden) layer:  <math>a_j = \sum_{i=0}^D w_{ji}^{(1)} x_i</math> and <math>z_j = h(a_j)</math>  second (output) layer:  <math>a_k = \sum_{j=0}^M w_{kj}^{(2)} z_j</math> and <math>y_k = h(a_k)</math>  =&gt; <math>y_k = \sigma \left( \sum_{j=0}^M w_{kj}^{(2)} h \left( \sum_{i=0}^D w_{ji}^{(1)} x_i \right) \right)</math></p> <p><b>three layer example</b>  first (hidden) layer:  <math>a_j = \sum_{i=0}^D w_{ji}^{(1)} x_i</math> and <math>z_j = h(a_j)</math>  second (hidden) layer:  <math>a_p = \sum_{j=0}^M w_{pj}^{(2)} z_j</math> and <math>z_p = h(a_p)</math>  third (output) layer:  <math>a_k = \sum_{p=0}^P w_{kp}^{(3)} z_p</math> and <math>y_k = h(a_k)</math></p> <p><b>Matrix notation:</b>  <math>a_{u_l} = W^{(l)} z_{u_{l-1}}</math>  <math>a_{u_l}</math>: <math>D_l</math>-sized vector of linearly combined units from previous vector:  <math>a_{i_l} = w_{i_l 1_{l-1}} z_{1_{l-1}} + w_{i_l 2_{l-1}} z_{2_{l-1}} + \dots + w_{i_l D_{l-1}} z_{D_{l-1}}</math>  <math>a_{D_l} = w_{D_l 1_{l-1}} z_{1_{l-1}} + w_{D_l 2_{l-1}} z_{2_{l-1}} + \dots + w_{D_l D_{l-1}} z_{D_{l-1}}</math>  <math>W^{(l)}</math>: a <math>D_l \times D_{l-1}</math> weight matrix  <math>z_{u_{l-1}}</math>: <math>D_{l-1}</math>-sized vector of hidden units</p> <p><b>Alternative indices:</b>  Each layer has hidden units <math>z_{l_{\text{out}}} = h(a_{l_{\text{out}}})</math> and a lin. comb. of input <math>z_{l_{\text{in}}}</math> with weights <math>w_{l_{\text{out}} l_{\text{in}}}^{(\text{layer})}</math>:  <math>a_{l_{\text{out}}} = \sum_{l_{\text{in}}=1}^{l_{\text{in}}} w_{l_{\text{out}} l_{\text{in}}}^{(\text{layer})} z_{l_{\text{in}}}</math></p>

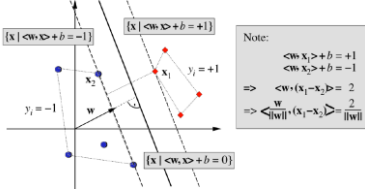
Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Conv. Neural Networks (CNN)</b> 	ML11.x	<p><b>Convolution:</b> filter <math>w: W^T + b</math></p>  <p>slide filter over image and compute (5x5x3 + bias) dot product at each location (5x5x3) filter weights are shared across activations of one activation map</p> <p>often implemented as cross-correlation (difference in mathematics)</p> <p><b>Nonlinearity:</b> nowadays usually ReLU</p> <p><b>Pooling:</b> (~subsampling) (not trained) makes representations smaller and more manageable; operates over each activation map independently; pooling filter size makes representation approximately invariant to small translations in the input nowadays usually max-pooling important to keep convolution cost low</p> <p><b>Regularization</b> through weight sharing</p>	<p>working with volumes (instead of vectors)</p> <p><b>hyperparameters:</b> <b>stride:</b> how fine grained you want to sample input layer <b>Layer output size:</b> <math>\left(\frac{N-F}{stride}\right) + 1</math> Layer input dim <math>N</math>, Filter dim <math>F</math></p> <p><b>valid convolution:</b> no padding and kernel required to be fully contained in image (but decreases output volume) <b>same convolution:</b> with padding to keep image size constant through successive convolutions</p> <p><b>zero padding:</b> with <math>\frac{F-1}{2}</math> (to preserve image size)</p> <p><b>Free parameters:</b> <math>\sum_{l=1}^L \dots</math></p> <p><b>nr of parameters per layer:</b> filter dimensions * output dimension: <math>[F \times F \times N_{layer\ l-1} + 1 (= bias)] \times N_{layer\ l}</math></p>	<p>unsupervised Learning</p> <p>Nonlinear Regression / Classification / Multiclass Classification (<math>K</math> output units and softmax activation to <math>y_k = \frac{\exp(a_k(x;w))}{\sum_j \exp(a_j(x;w))}</math>)</p> <p>• Image Classification / Semantic Segmentation / Retrieval / Captioning / ... • Pose estimation / view generation / ... • Playing games</p>	<p><b>training:</b> [feedforward &amp; backward pass] * nr. of gradient descent steps</p> <p><b>test/recall:</b> feedforward pass</p>	<p>+ tremendously successful in practical applications + shared parameters (built-in regularization)</p> <p>- expensive training - many hyperparameters / complex tuning - not robust to hyperparameter settings</p>	<p>Recurrent Neural Network (RNN)</p> <p>Multilayer Perceptron (MLP) fully connected version</p>
<b>Recurrent Neural Network (RNN)</b> 		<p><b>Allow for cycles in connectivity graph</b></p> <p>input <math>u(t)</math> <math>x(t) = f(\text{net}^x(t))</math> [nonlinearity <math>f</math>] <math>\text{net}^x(t) = Wx(t-1) + W^{\text{in}} u(t)</math> [recurrent weights <math>W</math>, input weights <math>W^{\text{in}}</math>]</p>  <p>output <math>y(t) = f^{\text{out}}(\text{net}^y(t))</math> [<math>f^{\text{out}}</math> output activation function] <math>\text{net}^y(t) = W^{\text{out}} x(t)</math></p> <p><b>Backpropagation Through Time (BPTT)</b> Idea: unfold network over time (results in one very deep network), then do backpropagation with weight sharing (at each step the same weights are used)</p>	<p><b>BPTT problems:</b> error gradient is propagated back and multiplied with weight <math>w</math> and derivative of unit's activity <math>f'(net)</math> <b>exploding gradient:</b> If <math> f'(net)w  &gt; 1 </math> for all <math>t</math>, exponential increase =&gt; oscillating weights <b>vanishing gradient:</b> If <math> f'(net)w  &lt; 1 </math> for all <math>t</math>, exponential decrease =&gt; slow convergence =&gt; cure: LSTM</p>	<p>unsupervised Learning</p> <p>Nonlinear Regression / (Multiclass) Classification (<math>K</math> output units and softmax activation to <math>y_k = \frac{\exp(a_k(x;w))}{\sum_j \exp(a_j(x;w))}</math>)</p> <p>• processing sequences • dynamical systems (rather than function mappings)</p>		<p>+ tremendously successful in practical applications + shared parameters (built-in regularization) + arbitrary complex</p> <p>- expensive to train - many hyperparameters / complex tuning - not robust to hyperparameter settings - vanishing gradient → LSTM - exploding gradient</p>	

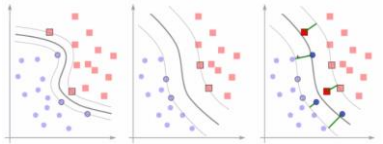
Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Classification and Regression trees (CART)</b>	ML8.x (regression)  ML8.21ff (classification)	<p>CART(<math>X, t, \text{max\_depth}, \text{min\_leaf}</math>)</p> <ol style="list-style-type: none"> <li>1: Check whether data should be split further; otherwise return leaf node</li> <li>2: Find best split value for each feature → Greedily minimize sum of squared errors in the two children</li> </ol> $x_{\text{split}} = \arg \min_x \left( \sum_{x_i \leq x} (t_i - \hat{y}_{\text{left}})^2 + \sum_{x_i > x} (t_i - \hat{y}_{\text{right}})^2 \right)$  <ol style="list-style-type: none"> <li>3: Choose best combination of split feature and split value (w.r.t. to squared error)</li> <li>4: Split data into left and right accordingly: (<math>X_l, t_l</math>) and (<math>X_r, t_r</math>)</li> <li>5: Save split feature and value, and pointer to two new subtrees to be</li> <li>6: built recursively: ( CART(<math>X_l, t_l, \text{max\_depth}-1, \text{min\_leaf}</math>), CART(<math>X_r, t_r, \text{max\_depth}-1, \text{min\_leaf}</math>) )</li> </ol>	<p>Input: data <math>X</math>, targets <math>t</math></p> <p><b>hyperparameters</b></p> <p>min_leaf: nr of samples (in leaf) max_depth: of the tree Total number of nodes Split model weak learner: axis-aligned splits → constant, other split models) Leaf/split model (majority vote; probability of data in the leaf) Split criterion Gini index; variance reduction; information gain</p> <p><b>entropy</b> “disorder/uncertainty” of Random Var. <math>V</math> with <math>K</math> possible outcomes <math>v_k</math> and distr <math>p(v_k)</math></p> $H(V) = - \sum_{k=1}^K p(v_k) \log_2 p(v_k)$ <p><b>information gain:</b> <math>V</math> has <math>N</math> data points with; split <math>s</math> into nodes <math>N_l</math> and <math>N_r</math>, RVs <math>V_l</math> and <math>V_r</math> with <math>p_l</math> and <math>p_r</math>:</p> $I = N \cdot H(V) - N_l \cdot H(V_l) - N_r \cdot H(V_r)$ <p>High-variance, low bias model (small changes of data may result in very different trees)</p> <p><b>Free parameters:</b> ...</p>	<p>supervised Learning unsupervised Learning</p> <p>Nonlinear Regression / (Multiclass) Classification</p> <ul style="list-style-type: none"> <li>• semi-supervised learning</li> <li>• density estimation</li> <li>• embedding learning</li> </ul>	<p><b>training:</b> ...</p> <p><b>test/recall:</b> ...</p> <p><b>storage during / after building:</b> <math>O(N)</math> / <math>O(ND)</math></p> <p><math>N</math> data points of dimensionality <math>D</math>, axis-aligned splits</p>	<p>+ easy to interpret + directly handle categorical features + scalable to large datasets (fast) + flexible framework with exchangeable components (split criterion, leaf model, type of split)</p> <p>– Tend to overfit – deterministic (i.e. not suitable for some ensemble methods)</p> <p>High-variance, low bias model (small changes of data may result in very different trees)</p> <p>random forests are said to be the best off-the-shelf model (for many applications)</p>	
<b>Random Forests</b>		<p>RF for CART</p> <ol style="list-style-type: none"> <li>1: For <math>b = 1</math> to <math>B</math></li> <li>2: Draw bootstrap sample <math>Z^*</math> of size <math>N</math> from training data</li> <li>3: Grow a random-forest tree <math>T_b</math> to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size <math>n_{\text{min}}</math> is reached:</li> <li>4: select <math>m</math> variables at random from <math>p</math> variables</li> <li>5: pick the best variable/split-point among the <math>m</math></li> <li>6: split the node into two daughter nodes</li> <li>7: Output the ensemble of trees <math>\{T_b\}_{1 \leq b \leq B}</math></li> </ol> <p>prediction for new point: Regression: <math>f^{\wedge}_{\text{rf}}(x) = 1/B \sum_{b=1:B} T_b(x)</math> Classification: Let <math>C^{\wedge}_b(x)</math> be the class pred. of the <math>b</math>-th random-forest tree. Then <math>C^{\wedge}_{\text{rf}}(x) = \text{majority vote } \{C^{\wedge}_b(x)\}_{1 \leq b \leq B}</math></p>	<p><b>hyperparameters</b></p> <p>CART-hyperparameters, additional: <math>m</math>: nr of (randomly chosen) variables at each split; <math>B</math>: Total number of trees</p> <p><b>Free parameters:</b> ...</p> <p>A1: individual models perform (reasonably) well A2: different models are not correlated</p>	<p>supervised Learning unsupervised Learning</p> <p>Nonlinear Regression / (Multiclass) Classification</p> <p>Motivation: Ensemble error depends on:</p> <ul style="list-style-type: none"> <li>• strength of indiv. models (→ we want indiv. good working models)</li> <li>• uncorrelatedness of models' errors (→ ... with high variance)</li> </ul>	<p><b>training:</b> ...</p> <p><b>test/recall:</b> ...</p> <p>Best case: balanced trees Fitting: <math>O(BmN \log N)</math> Prediction: <math>O(B \log N)</math> Worst case: splitting o one data point at a time Fitting: <math>O(BmN^2)</math> Prediction: <math>O(BN)</math></p> <p><b>storage during / after building:</b> <math>O(N)</math> / <math>O(ND)</math></p> <p><math>N</math> data points of dimensionality <math>D</math>, axis-aligned splits</p>	<p>+ robust to hyperparameter settings (vs. neural networks) + robust performance even for small datasets + bootstrapping: robust to outliers (and decorrelated trees)</p> <p>– Relatively weak performance for smooth functions without noise</p> <p>Out-of-bag error with little overhead: not every data point is used to t every single tree in fact, almost 37% are not used in each tree (see assignment) predict unused points for each tree to get unbiased estimated of the generalization error</p>	 <p>Decision Tree + (best of a fixed nr of) random splits = <b>ExtraTree</b> Extra Tree + Bagging = <b>ExtraTrees</b> Decision Trees + Bagging = <b>Bagged trees</b> Decision Trees + best split using a random subset of <math>m \leq D</math> features + Bagging = <b>Random forest</b></p>

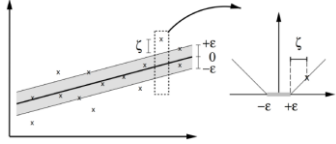
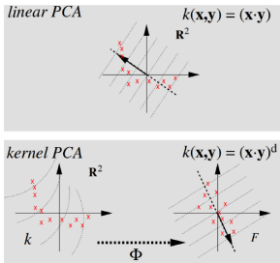
Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Boosting / AdaBoost</b>	ML6.36  ML9.x	<p>Ensemble technique, submodels are trained sequentially, with the m-th submodel trained to fix the mistakes of the previous (first m-1) submodels</p> <p>AdaBoost Binary classification class, labels +1 and -1 Final model <math>G(x)</math> combines individual submodels <math>G_1, \dots, G_M</math>, through a weighted majority vote with weights <math>\alpha_1, \dots, \alpha_M</math>: <math>G(x) = \text{sign}(\sum_{m=1}^M \alpha_m G_m(x))</math></p> <p>with weights <math>\alpha_m(\text{err}_m) = \log(\frac{1-\text{err}_m}{\text{err}_m})</math>  <math>(\alpha_m(0.5) = 0, 0 \leq \text{err}_m \leq 1 \rightarrow \alpha_m \rightarrow \pm 10)</math></p> <p>weighted training error rate: <math>\text{err}_m = \frac{\sum_{i=1}^N w_i^m I(y_i \neq G(x_i))}{\sum_{i=1}^N w_i^m}</math></p> <p>with <math>w_i^{m+1} = w_i^m \cdot \exp(\alpha_m I(y_i \neq G(x_i)))</math> being adapted weights for each iteration <math>m</math> and indicator function <math>I</math>:  <math>I(a) = \begin{cases} 1, &amp; \text{if } a = \text{true} \\ 0, &amp; \text{else} \end{cases} \Rightarrow \begin{cases} w_i^m \\ w_i^m \cdot \exp(\alpha_m) \end{cases} = w_i^{m+1}</math></p> <p>AdaBoost algorithm (weights <math>w_i^m</math> don't need to be stored, may be overwritten)  1: init weights <math>w_i^0 = \frac{1}{N}, i = 1, 2, \dots, N</math>  2: for <math>m = 1</math> to <math>M</math>:  3: fit classifier <math>G_m(x)</math> to <math>w_i^{m-1}</math>-weighted train data  4: compute <math>\text{err}_m = \frac{\sum_{i=1}^N w_i^m I(y_i \neq G(x_i))}{\sum_{i=1}^N w_i^m}</math>  5: compute <math>\alpha_m(\text{err}_m) = \log(\frac{1-\text{err}_m}{\text{err}_m})</math>  6: set <math>w_i^{m+1} = w_i^m \cdot \exp(\alpha_m I(y_i \neq G(x_i)))</math>  7: output <math>G(x) = \text{sign}(\sum_{m=1}^M \alpha_m G_m(x))</math></p> 	<p><b>hyperparameters</b> CART-hyperparameters, additional:  <math>m</math>: nr of (randomly choosen) variables at each split;  <math>B</math>: Total number of trees</p> <hr/> <p><b>Free parameters:</b> ...</p> <hr/>			<p>+ good generalization capabilities on low noise data</p> <p>- if data is very noisy, it may overfit badly</p> <p>successfully used in practice</p> <p>boosting (like RF) said to be the best off-the-shelf model (for many applications)</p>	Ensemble technique like bagging, but models are not trained independently
<b>Forward Stagewise Additive Modelling (FSAM)</b>	ML9.25ff	<p>loss function <math>L(\cdot, \cdot)</math>: quantify error between our model predictions <math>f(x_i)</math> and targets <math>y_i</math>  E.g., squared error: <math>L(y_i; f(x_i)) = 1/2(y_i - f(x_i))^2</math></p> <p>way to construct a model to greedily minimize a loss function  - E.g., a routine to find the parameters of some model <math>b</math>:  <math>m = \text{argmin}_b \sum_{i=1}^N L(y_i; b(x_i))</math></p>				+	can be shown to recover the AdaBoost algorithm.

Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Gradient Tree Boosting</b>	ML9.30ff	<p>minimize this: <math>L(f) = \sum_{i=1}^N L(y_i, f(x_i))</math> w.r.t. <math>f</math> (constraint A2)</p> <p>One can see minimization of training error as gradient descent in the N-dimensional space of “parameters” that describe the values of <math>f</math> at the N data points <math>x_i</math></p> <p><math>\hat{f} = \operatorname{argmin}_f (L(f))</math> with <math>f = [f(x_1), \dots, f(x_N)]^T</math></p> <p>Gradient Tree Boosting Algorithm</p> <ol style="list-style-type: none"> <li>1: init weights <math>f_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^N L(y_i, \gamma)</math></li> <li>2: for <math>m = 1</math> to <math>M</math>:</li> <li>3:   for <math>i = 1, 2, \dots, N</math>              compute <math>r_{im} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}</math></li> <li>4:   fit a regression tree to targets <math>r_{im}</math> giving terminal regions <math>R_{jm}, j = 1, 2, \dots, J_m</math></li> <li>5:   for <math>i = 1, 2, \dots, N</math> compute              <math>\gamma_{jm} = \operatorname{argmin}_\gamma \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma)</math></li> <li>6:   update <math>f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})</math></li> <li>7:   output <math>\hat{f}(x) = f_M(x)</math></li> </ol>	<hr/> <p><b>Free parameters:</b> ...</p> <hr/> <p>A1: Loss function is differentiable  A2: <math>f</math> comes from a certain model class (e.g., a sum of trees)</p>	<p>supervised Learning unsupervised Learning</p> <p>Nonlinear Regression / (Multiclass) Classification</p> <hr/> <p>Motivation:</p> <ul style="list-style-type: none"> <li>• computationally effective boosting methods only for exp/squared error loss</li> <li>• With gradient boosting: any differentiable loss function works</li> </ul>	<p><b>training:</b>  <math>O_{best}(BmN \log N)</math>  <math>O_{worst}(BmN^2)</math></p> <p><b>test/recall:</b>  <math>O_{best}(B \log N)</math>  <math>O_{worst}(BN)</math></p> <p><b>Best case:</b>  balanced trees  <b>Worst case:</b>  splitting one data point at a time</p> <p><b>storage during / after building:</b>  <math>O(N)</math> / <math>O(ND)</math></p> <p><math>N</math> data points of dimensionality <math>D</math>, axis-aligned splits</p>	<p>+ robust to hyperparameter settings (vs. neural networks)  + robust performance even for small datasets  + bootstrapping: robust to outliers (and decorrelated trees)</p> <p>– Relatively weak performance for smooth functions without noise</p> <p>Out-of-bag error with little overhead:  not every data point is used to train every single tree  in fact, almost 37% are not used in each tree (see assignment)  predict unused points for each tree to get unbiased estimated of the generalization error</p>	



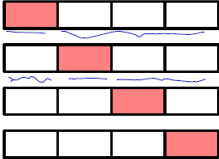
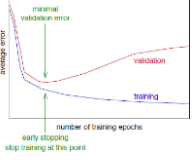
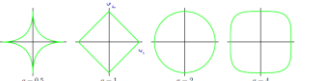
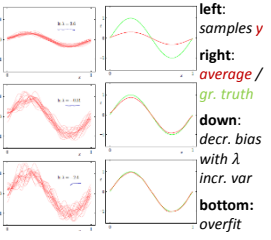
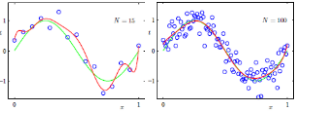
Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Large Margin classifier</b>  Optimal hyperplane estimator	ML13.x	 <p><b>Large Margin Classifier</b>            optimize <math>w, b</math> to get the largest possible smallest distance between data points and decision hyperplane <math>x</math>:  <math>\arg \max_{w \in \mathcal{H}, b \in \mathbb{R}} \min(\ x - x_i\ )</math> for <math>i = 1, \dots, m</math>            with <math>x \in \mathcal{H}</math> and <math>\langle w, x \rangle + b = 0</math> (decision function)            This is solved by the shortest normal vector <math>w</math> (that still gives correct classification results)  <math>\arg \min_{w \in \mathcal{H}, b \in \mathbb{R}} \tau(w) = \frac{1}{2} \ w\ ^2</math> (<b>objective function</b> <math>\tau</math>)            subject to <b>inequality constraint</b> <math>y_i(\langle w, x_i \rangle + b) \geq 1</math> with <math>i = 1, \dots, m</math> (<math>\geq 1</math> fixes the scaling of <math>w</math>, data should be a bit away from decision plane, 1 could be any other number)</p> <p><b>Lagrangian:</b> first part should be small, second part large  <math display="block">L(w, b, \alpha) = \frac{1}{2} \ w\ ^2 - \sum_{i=1}^m \alpha_i y_i (\langle w, x_i \rangle + b) - 1</math>            approach:            use partial derivatives, minimize Lagrangian <math>L</math> w.r.t. primal variables <math>w</math> and <math>b</math> and maximize w.r.t. dual variables <math>\alpha_i</math> (effectively finds solution in a saddle point)  <math>\frac{\partial}{\partial b} L(w, b, \alpha) = 0 \rightarrow \sum_{i=1}^m \alpha_i y_i = 0</math>  <math>\frac{\partial}{\partial w} L(w, b, \alpha) = 0 \rightarrow \sum_{i=1}^m \alpha_i y_i x_i = w</math>            eliminate primal variables by substitution (of <math>w</math> and <math>b</math>)            dual optimization problem of vector <math>\alpha</math> (quadratic program):  <math display="block">\arg \max_{\alpha \in \mathbb{R}^m} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle</math>            subject to <math>\alpha_i \geq 0</math> for all <math>i = 1, \dots, m</math> and <math>\sum_{i=1}^m \alpha_i y_i = 0</math>            hyperplane decision function “instance based learning”  <math display="block">f(x) = \text{sign} \left( \sum_{i=1}^m \alpha_i y_i \langle x, x_i \rangle + b \right)</math>            just working on dot products of data.            Vector <math>\alpha</math> may be sparse: a lot of <math>\alpha_i = 0</math>, corresponding <math>x_i</math> do not contribute to decision hyperplane. Otherwise they are called Support Vectors (SV).</p>		supervised Learning  Classification <i>linearly seperable case</i>  <hr/> Motivation: <ul style="list-style-type: none"> <li>basis for SVM</li> </ul>		+ just working on dot products, no assumptions about data vs + maximizing the margin also restricts the function class (simple/smooth functions are preferred)  - ...	instance based learning:  <b>k-Nearest Neighbor</b> (drawback: save all data points, here: only ones with “active” lagrange multipliers $\alpha_i$ that actually influence the decision hyperplane; corresponding training data points are called support vectors)

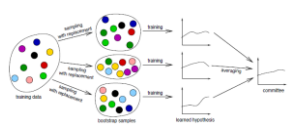
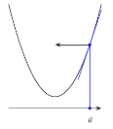
Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
<b>Support Vector Machines (SVM)</b>	ML13.x ML14.x  Bishop 7	<p>See <b>Large Margin classifier</b> (generalization by kernel trick): optimize <math>w, b</math> to get the largest possible smallest distance between data points and decision hyperplane <math>x</math> in non-linear space</p> $\operatorname{argmax}_{\alpha \in \mathbb{R}^m} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j k(x_i, x_j)$ <p>subject to <math>\alpha_i \geq 0</math> for all <math>i = 1, \dots, m</math> and <math>\sum_{i=1}^m \alpha_i y_i = 0</math></p> <p><b>Building Blocks:</b>  <b>similarity measure <math>k</math> / “kernel”</b>  for: arbitrary input space, e.g. set <math>x \in \mathcal{X}</math> (previously we had: <math>x \in \mathbb{R}^N</math> vector space with defined similarity)  (1) delivers a real number describing similarity <math>k</math> between two patterns <math>k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^N, (x, x') \rightarrow k(x, x')</math>  (2) is symmetric <math>k(x, x') = k(x', x)</math>  define from dot product of a feature space <math>\mathcal{H} \in \mathbb{R}^N</math>:  <math>k(x', x) := \langle x, x' \rangle = \langle \phi(x), \phi(x') \rangle</math> where <math>\phi: \mathcal{X} \rightarrow \mathcal{H}</math> maps <b>input space <math>\mathcal{X}</math> to feature space <math>\mathcal{H}</math></b></p> <p><b>Large margin classifier</b>  optimize <math>w, b</math> to get the largest possible smallest distance between data points and decision hyperplane <math>x</math></p> <p>linear case:  <math>f(x) = \operatorname{sign}(\sum_{i=1}^m \alpha_i y_i \langle x_i, x \rangle + b)</math></p> <p>with kernel trick:  <math>f(x) = \operatorname{sign}(\sum_{i=1}^m \alpha_i y_i \langle \phi(x_i), \phi(x) \rangle + b)</math>  <math>= \operatorname{sign}(\sum_{i=1}^m \alpha_i y_i k(x, x_i) + b)</math></p>	<p><b>hyperparameters:</b></p> <p>kernel parameters:  polynomial kernel (param <math>d</math>)  <math>k(x', x) = \langle x, x' \rangle^d</math>  Gaussian radial basis functions (param <math>\sigma</math>)  <math>k(x', x) = \exp\left(-\frac{1}{\sigma^2} \ x, x'\ ^2\right)</math>  sigmoid kernel (param <math>k &gt; 0</math> and <math>\theta &lt; 0</math>)  <math>k(x', x) = \tanh(k \langle x, x' \rangle + \theta)</math></p> <hr/> <p><b>Free parameters:</b> ...</p> <hr/> <p>A1: All computations can be formulated in dot product space</p> <p><b>VC bound</b>  VC dimension <math>h</math> of function class: number of data points <math>m</math> you can scatter in feature space <math>\mathbb{R}^N</math> (proportional to capacity)  <math>R[f] \leq R_{\text{emp}}[f] + \sqrt{\frac{1}{m} \left( h \left( \log \frac{2m}{h} + 1 \right) - \log \frac{\delta}{4} \right)}</math>  risk is bounded by training error + capacity</p> <p>risk w.r.t. to observed data (“empirical risk” / “average training error”):  <math>R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{2}  f(x_i) - y_i  \right)</math>  risk w.r.t. to underlying data distribution:  <math>R[f] = \int \frac{1}{2}  f(x) - y  dP(x, y)</math></p>	<p>supervised Learning</p> <p>Classification</p> <hr/> <p>Motivation:</p> <ul style="list-style-type: none"> <li>make use of nonlinear mappings but don't compute dot product in feature space explicitly</li> </ul>	<p><b>training:</b> ...</p> <p><b>test/recall:</b> ...</p>	<p>+ good in large feature space  + very robust performance, even when used as black box  + very few hyperparameters  + inspecting selected support vectors may help to understand the problem  + SVM formulation can be extended to non-separable cases (in feature space)  → soft-margin SVM  + SVM regression and outlier detection (“one-class SVM”) easily obtained  + dealing with unbalanced classes by weighting influence of misclassifications separately per class</p> <p>– kernel computations expensive for many training data points  → chunking methods avoid calculating full kernel matrix</p> <p>– bad inspection properties (vs. linear methods)</p> <p>– outperformed by Deep NN (at least if dataset is huge)</p>	<p>Relevance SVM (Bishop 7.2)</p> <p>soft-margin SVM  slack variables <math>\xi_i \geq 0</math> per data point  <math>y_i \langle w, x_i \rangle + b \geq 1 - \xi_i</math>  <math>\operatorname{argmin}_{w \in \mathcal{H}, b \in \mathbb{R}} \tau(w) = \frac{1}{2} \ w\ ^2 + C \sum_{i=1}^m \xi_i</math>  hyperparameter <math>C</math>: trade-off between enlarging the margin and minimizing training error</p>
<b>Soft-margin SVM</b>	ML14.x	<p><math>\operatorname{argmax}_{\alpha \in \mathbb{R}^m} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j k(x_i, x_j)</math></p> <p>subject to <math>0 \leq \alpha_i \leq C</math> for all <math>i = 1, \dots, m</math> and <math>\sum_{i=1}^m \alpha_i y_i = 0</math></p> <p>slack variables <math>\xi_i \geq 0</math> per data point  <math>y_i \langle w, x_i \rangle + b \geq 1 - \xi_i</math>  <math>\operatorname{argmin}_{w \in \mathcal{H}, b \in \mathbb{R}} \tau(w) = \frac{1}{2} \ w\ ^2 + C \sum_{i=1}^m \xi_i</math></p> <p>allows for slight violation (penalty must be determined at training time with regularization parameter <math>C</math> resp. <math>\nu</math>)</p>  <p>1 – hard margin (complex hyperplane but no errors)  2 – without original SVM  3 – soft margin SVM (potentially smoother hyperplane)</p>	<p>same as SVM</p> <p><math>b</math> (bias scalar): can be computed exploiting that for all SVs <math>x_i</math> with <math>\alpha_i &lt; C \rightarrow \xi_i = 0</math>.  Intuition: shifts the decision hyperplane such, that SVs with zero slack lie on <math>\pm 1</math> lines.</p> <p><b>hyperparameters:</b>  <math>C</math>: trade-off between enlarging the margin and minimizing training error  <b>small <math>C</math></b> <math>\rightarrow</math> may need more SV to define hyperplane (since they have small influence)  <b>large <math>C</math></b> <math>\rightarrow</math> fewer SV but with large weights</p> <p><math>\nu</math>: alternative parametrization with <math>0 &lt; \nu \leq 1</math>:  bounds the fraction of training patterns which will become SVs and those which will have non-zero slacks.</p> <hr/> <p><b>Free parameters:</b> ...</p> <hr/>	<p>supervised Learning</p> <p>Classification</p> <hr/> <p>Motivation:</p> <ul style="list-style-type: none"> <li>noisy data</li> <li>misclassified data</li> </ul>	<p><b>training:</b> ...</p> <p><b>test/recall:</b> ...</p>	<p>same as SVM</p> <p>Regularization:  + allows solutions with function class of lower capacity (better generalization)  + enlarged margin  + reduce <math>R[f]</math></p> <p>– outperformed by Deep NN (at least if dataset is huge)</p>	

Method	slides	Idea	Characteristics / Parameters / Assumptions	Usage / Motivation	Runtime	Pro + and Contra –	special cases / improvements / related methods
Support Vector Regression (SVR)	ML14.27ff	<p>a support vector regression model can estimate continuous non-linear functions <math>\mathbb{R}^N \rightarrow \mathbb{R}</math>  Idea: <math>\varepsilon</math>-tube around function</p> <p>Key ingredients:  Kernel trick  Use the <math>\varepsilon</math>-insensitive loss function (corresponds to regularizing the solution similar to a large margin).  Introduce (two types) of slack variables <math>\xi_i</math>, which allow for violations of the <math>\varepsilon</math>-tube. Limit influence of outliers by introducing a constant C.  Choose constants <math>C, \varepsilon \geq 0</math> a priori  Formulation via Lagrange multipliers, solve quadratic problem.</p> 		<p>supervised Learning</p> <p>(non-linear) Regression</p> <hr/>			<a href="http://www.svms.org/regression/SmSc98.pdf">http://www.svms.org/regression/SmSc98.pdf</a>
Kernel Principal Component Analysis (kPCA)	ML14.32ff	<p>Idea: Formulate all calculations of linear PCA (eigenvalue problem!) via dot products of the input space X. obtain a non-linear version of linear PCA by mappings - from input space X to a high-dimensional feature space H!  use the kernel trick, thus compute everything in input space X!</p> $f_n(x) = \sum_{i=1}^m \alpha_i^n k(x_i, x)$ <p>where <math>\alpha_i^n</math> are (up to a normalizing constant) the components of the <math>n^{\text{th}}</math> eigenvector of the Gram Matrix <math>K_{ij} := (k(x_i, x_j))</math> with <math>\dim(K): M \times M</math> and <math>n = 1, \dots, m</math> (as many as data points)</p>  <p>linear PCA <math>k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})</math></p> <p>kernel PCA <math>k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^d</math></p> <p>input space <math>\mathcal{X}</math> to feature space <math>\mathcal{H}</math></p>	<p><b>hyperparameter:</b>  <math>M</math>: nr. of output dim</p> <hr/> <p><b>Free parameters:</b> ...</p> <hr/> <p>A1: ...</p>	<p>non-linear feature extractor</p> <hr/> <p>pre-process for non-kernel methods:</p> <ul style="list-style-type: none"> <li>• data compression</li> <li>• dimensionality reduction</li> </ul> <p>introspection:</p> <ul style="list-style-type: none"> <li>• into SVM solution: posthoc vis. of variance isolines (of feature func. in <math>\mathcal{X}</math> using SVM hyperparams)</li> <li>• why a classification problem is hard ([1] intrinsically complex/high dim or [2] simple but noisy?)</li> </ul>	<p><b>training:</b>  <math>O(N^3)</math>  (<math>N, M</math> or <math>m</math>: nr of data points)</p> <p>or iterative methods that scale better for high dim.</p> <p><b>test/recall:</b>  ...</p>	<p>+ allows non-linear feature reduction</p> <p>- computation of Gram matrix expensive for large nr of patterns</p>	

## 2. General Machine Learning Concepts

Concept	slides	Idea	Characteristics / Parameters / Assumptions / Examples	Usage / Motivation		special cases / related methods / improvements
Approximat.-Generaliz. Tradeoff	ML6.8	<p><b>Approximation-Generalization Tradeoff</b>  more complex &lt;&gt; more flexibility to approx..  less complex &lt;&gt; more likely to generalize</p> <p>squared distance error:  <math>E(w) = \frac{1}{2} \sum_{n=1}^N (y(x_n, w) - t)^2</math>  root-mean-square (RMSE) normalized error  <math>E_{RMS}(w) = \sqrt{2E(w^*)/N}</math>  N: allows to compare different sizes of data  sqrt: makes sure E_RMS is measured in same scale as t</p>				
Loss function	ML6.11	<p>expected loss (regression with squared error)  <math>\mathbb{E}[L] = \int (y(x) - h(x))^2 p(x) dx + \int (h(x) - t)^2 p(x, t) dx</math>  first term: how much does the prediction function ("hypothesis") <math>y(x)</math> differ from the ideal (and theoretical) prediction function <math>h(x)</math>  second term: how much does the ideal prediction function differ from the real target labels <math>t</math>. This is the noise inherent to the data.</p>		Supervised Learning		
Bias-Variance	ML6.12ff	<p><math>\mathbb{E}_{\mathcal{D}} [(y(x, \mathcal{D}) - h(x))^2] = (\bar{y}(x) - h(x))^2 + \mathbb{E}_{\mathcal{D}} [(y(x, \mathcal{D}) - \bar{y}(x))^2]</math>  average performance of <i>prediction function</i> ("hypothesis") <math>y(x, \mathcal{D})</math> over multiple <i>datasets</i> <math>\mathcal{D}</math>; extended and reformulated with <i>average prediction function</i> <math>\bar{y}(x) = \mathbb{E}_{\mathcal{D}}[y(x, \mathcal{D})]</math>.</p> <p>(<b>bias</b>)<sup>2</sup>: how much differs the (<i>average</i>) <i>hypothesis</i> (over all <i>datasets</i>) <math>\bar{y}(x)</math> from the <i>ideal prediction function</i> <math>h(x)</math>  <b>variance</b>: difference of a single hypothesis to the average hypothesis; how much does the model vary with particular data <math>\mathcal{D}</math></p> <p><b>Bias-Variance Tradeoff</b> ~ Approximation-Generalization Tradeoff  less complex &lt;&gt; more likely to generalize (<b>high bias, low var</b>)  more complex &lt;&gt; more flexibility to approx.. (<b>low bias, high var</b>)  Balancing bias and variance gives optimal predictive capability.</p>	<p>left: <b>high bias, low/no var</b>  right: <b>low bias, high var</b></p>	Supervised Learning		<p><b>High bias-Low Variance:</b></p> <p><b>Low bias-High Variance:</b>  Trees</p>
Validation Set	ML6.17	<p>Require: data <math>\mathcal{D}</math>  1: split data to disjoint subsets: training set <math>\mathcal{D}_{train}</math>, test set <math>\mathcal{D}_{test}</math>  2: apply training only on the training set <math>\mathcal{D}_{train}</math>  3: apply testing of learned functions on the independent test/validation set <math>\mathcal{D}_{test}</math></p>		Supervised Learning	<p>+ independent set allows spotting overfitting  - only subset of avail. data is used for training/testing  → Cross-Validation</p>	<p><b>Related methods:</b>  Cross-Validation</p>

Concept	slides	Idea	Characteristics / Parameters / Assumptions / Examples	Usage / Motivation		special cases / related methods / improvements
Cross Validation (k-fold)	ML6.x	Require: data $\mathcal{D}$ , parameter $k$ : $2 \leq k \leq N$ ( $N$ : nr of data points) 1: split data set into $k$ disjoint subsets of equal size: $\mathcal{D}_1, \dots, \mathcal{D}_k$ 2: for $i = 1$ to $k$ do 3: train model on set $\mathcal{D}_1 \cup \dots \cup \mathcal{D}_{i-1} \cup \mathcal{D}_{i+1} \cup \dots \cup \mathcal{D}_k$ (all but $\mathcal{D}_i$ ) 4: calculate average test error $e_i$ on $\mathcal{D}_i$ 5: end for 6: return $\frac{1}{k} \sum_{i=1}^k e_i$ (average validation error)	<b>hyperparameter:</b> $k$ : number of segments error function $e_i$ , e.g. RMSE: $e = \sqrt{2E(w^*)/N} = \sqrt{2(y-t)^2}$ 	Supervised Learning  How to spot overfitting?	+ model is learned on $\frac{k-1}{k} N$ data points and evaluated on all - model has to be learned $k$ times	<b>Special Case:</b> leave-one-out Cross-Validation for $k = N$ if data is particularly scarce
Regularization	ML6.x	penalties (e.g. for large param values) preprocessing (remove outliers, in general filtering) more data (if errors are due to noise, increase sample size)				
Regularization Early stopping	ML6.25	stop learning when error on validation set has reached its minimum requires perpetual observation of val. error			+ -	
Regularization shrinkage methods	ML6.26	Extend error function with penalty term $E_W(w)$ $E(w) = E_D(w) + E_W(w)$ (e.g. for simple linear model $y = w^T x$ ) general form: sum of squares with Lq regularizer $E(w) = \frac{1}{2} \sum_{i=1}^N (t - w^T \phi(x_i))^2 + \lambda \frac{1}{2} \sum_{j=1}^M  w_j ^q$	<b>hyperparameters:</b> $\lambda$ : strength of regularization regularizer params (e.g. $q$ )  case $q = 1$ : L1 regularizer (Lasso: gives sparse solutions) case $q = 2$ : L2 regularizer (Ridge Regression / Tikhonov reg)		+ will enforce smaller weights - determine $\lambda$ (trial and error? k-fold cross val?) 	<b>Lasso</b> $E_W$ L1 regularization: gives sparse solutions <b>Ridge Regression / Tikhonov regularization</b> $E_D$ sum-of-squares with $E_W$ L2 regularization $E_W(w) = \lambda \frac{1}{2} w^T w$ allows analytic sol. $w = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T t$
Regularization more data	ML6.31	get more training data; if not possible directly, think about related sources of similar data				
Regularization Filtering	ML6.32	Reduce training data to a subset that is (ideally) free of noise and contains all important patterns for learning the task. Techniques: oversampling, subsampling, outlier rejection, jittering	Assume: Train data contains subset that doesn't contribute/distorts training (e.g. outliers)	unbalanced data in classification		
Regularization Feature Selection	ML6.32	feature extraction/selection: use more/less/other input features  less features: reduce overfitting by avoiding pseudo relationships more/other features: improve generalization if features are related to desired output non-linear transformations:			- must be chosen problem specific	<b>Related methods:</b> (semi-)automated feature selection  dimensionality reduction PCA, ICA, mutual information

Concept	slides	Idea	Characteristics / Parameters / Assumptions / Examples	Usage / Motivation		special cases / related methods / improvements
Regularization  Bootstrap		Robustness towards outliers Decorrelation of the trees in the ensemble Out-of-bag error: not every data point is used to train every single tree in fact, almost 37% are not used in each tree (see assignment) predict unused points for each tree to get unbiased estimated of the generalization error			+ heavily used in practice - ...  “bagged” estimator has <b>lower variance / higher bias</b> than the individual models it “bags”	
Regularization  Bagging (Bootstrap AGGregation)	ML6.35	train several models on bootstrap samples of the training data (data is drawn randomly with replacements) some patterns may occur twice or more, others don't occur at all average the output of all trained models  single members of the committee might produce a higher test-set error; however in general the diversity of the committee compensates for this effect and therefore the committee error improves over the error of the individuals			+ heavily used in practice - ...  “bagged” estimator has <b>lower variance / higher bias</b> than the individual models it “bags”	Tree based methods
Regularization  Boosting	ML6.36	ensemble technique like bagging, but models are not trained independently second model is learned on the training data that are not well learned by the first model third model is learned on the training data that are not well learned by the first and second model ... step by step, we get better committees on the training set boosting is successfully used in practice good generalisation capabilities on low-noise data			+	
Optimization  Gradient descent	ML7.25ff	Require: mathematical function $f$ , learning rate $\epsilon > 0$ Ensure: returned vector $u$ is close to a local minimum of $f$ 1: choose an initial point $u$ 2: while $\ grad f(u)\ $ not close to 0 do 3: $u \leftarrow u - \epsilon \cdot grad f(u)$ 4: end while 5: return $u$		Iterative Optimization  for many interesting minimize problems $argmin_u f(u)$ , no closed form solution (even though $\frac{\partial f}{\partial u_i} = 0$ holds at (local) solution points!) => iterative methods	+ ... - how to choose initial $u$ - how to choose $\epsilon$ - does this algorithm really converge?	