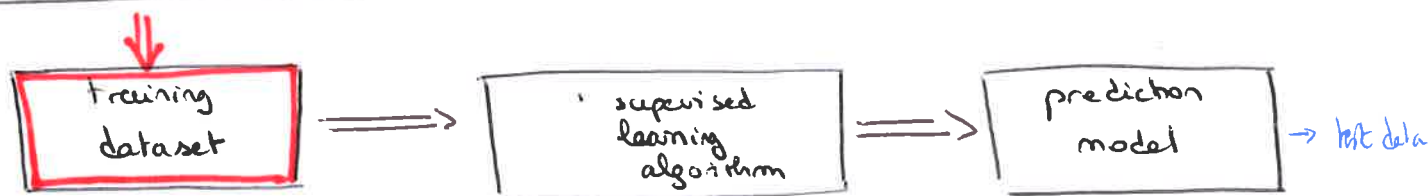


machine learning framework:



set of m input data:
 $T = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$
 with $x^{(i)} \in \mathbb{R}^m$ for $i = 1 \dots m$
 is the feature vector
 &
 $y^{(i)} \in \{1 \dots c\}$ for $i = 1 \dots m$
 is a class label

eg: SVM
 ANN
 LDA, ...

$h : \mathbb{R}^m \mapsto \{1 \dots c\}$
 $x \mapsto y$
 discrimination function trying to "mimic" the input-label relationship.

Problem: training dataset is of difficult / expensive to obtain

Goal: take advantage of unlabeled feature data set available:

$\{x_u^{(1)}, \dots, x_u^{(k)}\}, x_u^{(j)} \in \mathbb{R}^n$ for $j = 1 \dots k$.

eg: see semi-supervised learning [Nigam 2000]
 transfer learning [Thrun 1996, Caruana 1997, Ando 2005].

Proposal: improve the training dataset ^{Multi-task learning} using unlabeled data.
 feature set

- ① Use sparse coding to construct a succinct higher-level features representation using unlabeled data (only) ^{"abstract"}
- ② Apply this representation on labeled dataset and use it for supervised classification.

No assumption on the unlabeled dataset, in particular it doesn't assume that they have the same class labels (eg: task = Rhino vs Elephant, then can use unlabeled img of houses, beaches, horses, pigs... to try to improve discrimination).
 nor that feature data follow the same generative distribution.

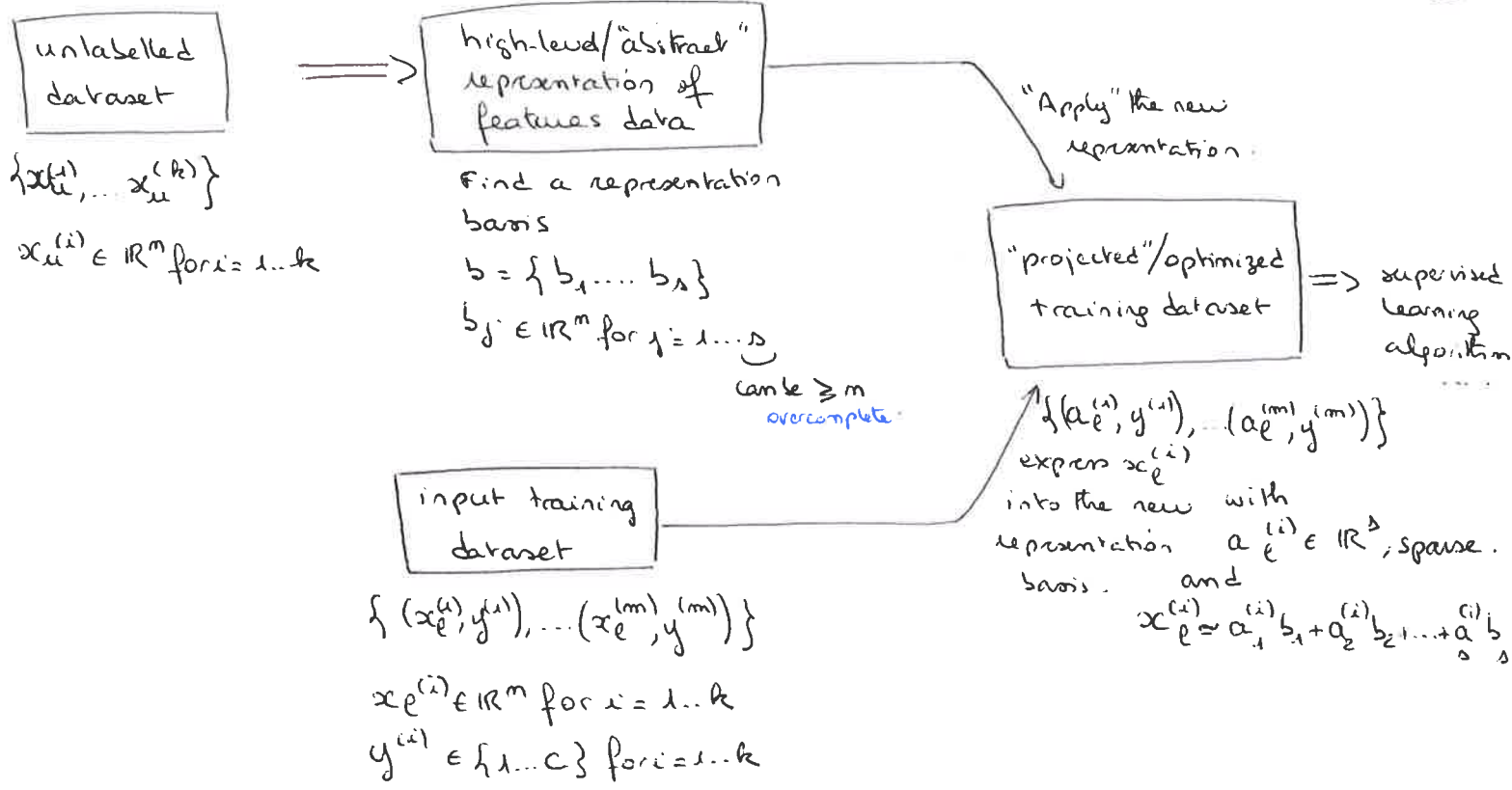
But it should be relevant eg. same modality, type of data.

Eg: images of Rhinos vs images of Elephants

can images of other scene (eg beach, house, horses, pigs, ...) improve classification of Rhinos vs Elephant?

⇒ help ^{detect} recognize / ^{structures} visual patterns in images (eg edges) ⇒ used for classification task.
 discover correlation between features
 salient features they share

Unsupervised feature construction:



1 Construct a representation basis for feature.

Solve:

(*)
$$\min_{b, a} \sum_{i=1}^k \|x_u^{(i)} - \sum_{j=1}^s a_j^{(i)} b_j\|_2^2 + \lambda \|a\|_1$$

such that $\|b_j\|_2 \leq 1, \forall j \in 1 \dots s$

constraints

good reconstruction of x_u as a weighted linear combination of b_j - loss function

L_1 - regularisation term to promote sparsity (\Rightarrow most elements $b=0$) (using another norm lead to lower performances).

\sim sparse coding algorithm

\Rightarrow come out with a set of basis vector $b = \{b_1, \dots, b_s\}$, $b_j \in \mathbb{R}^m$ for $j=1 \dots s$ (can $s > m$ overcomplete)

activations $a = \{a^{(1)}, \dots, a^{(k)}\}$, $a^{(i)} \in \mathbb{R}^s$ for $i=1 \dots k$

such that for $i=1 \dots k$:

$$x_u^{(i)} = a_1^{(i)} b_1 + a_2^{(i)} b_2 + \dots + a_j^{(i)} b_j + \underbrace{a_{j+1}^{(i)}}_0 b_{j+1} + \dots + \underbrace{a_s^{(i)}}_0 b_s \quad (\text{see Fig 263})$$

$x_u^{(i)}$ is a sparse weighted combination of b_j - only a few basis elements are used to reconstruct any input

(*) Solved iteratively by alternatingly optimizing over $a \perp b$.

② Express each training input $x_e^{(i)}$, $i=1 \dots m$, in the new feature representation ③ representation basis.

Find activations $a(x_e^{(i)}) \in \mathbb{R}^d$ for $i=1 \dots m$

Solve:

$$a(x_e^{(i)}) = \underset{a^{(i)}}{\operatorname{argmin}} \left\| x_e^{(i)} - \sum_{j=1}^d a_j^{(i)} b_j \right\|_2^2 + \beta \|a^{(i)}\|_1$$

note: no constraint, b is fixed, already optimized in step ①.

$$x_e^{(i)} = \begin{bmatrix} x_{e,1}^{(i)} \\ x_{e,2}^{(i)} \\ \vdots \\ x_{e,m}^{(i)} \end{bmatrix} \implies a(x_e^{(i)}) = \begin{bmatrix} a_1^{(i)} \\ a_2^{(i)} \\ \vdots \\ a_d^{(i)} \end{bmatrix} = 0$$

Sparse representation

Comparison to PCA

also used to model "higher-level" structure of the features dataset

Aim: identify a low-dimensional subspace of maximal variation defined by the $k \leq m$

Principal Component b_1, b_2, \dots, b_k solution of

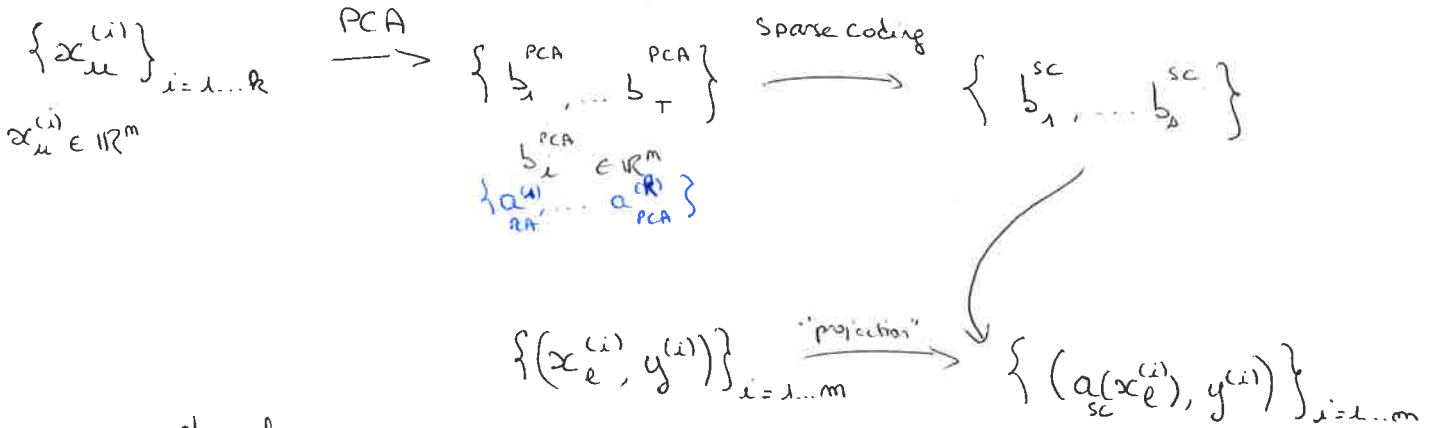
$$\underset{b, a}{\operatorname{argmin}} \sum_{i=1}^m \left\| x_e^{(i)} - \sum_{j=1}^k a_j^{(i)} b_j \right\|_2^2$$

such that b_1, b_2, \dots, b_k are orthogonal.

Advantage: finding $(a_j^{(i)})_{j=1 \dots k}$ corresponding to input $x^{(i)}$ is easy: $a_j(x^{(i)}) = b_j^T x^{(i)}$

- Limits:
- linear feature extraction (linear function of the inputs x)
 - b_j are assumed to be orthogonal
 - $k \leq m$

Experiments



compare classification learning performances using the following database for training: \rightarrow SVM, GDA

- $\{(a_{SC}(x_e^{(i)}), y)\}_{i=1 \dots m}$ ie unsupervised feature construct using sparse coding
- $\{(x_e^{(i)}, y)\}_{i=1 \dots m}$ ie raw features
- $\{(a_{PCA}(x_e^{(i)}), y)\}_{i=1 \dots m}$ ie PCA features or $\{([x_e^{(i)}, a_{PCA}(x_e^{(i)})], y)\}_{i=1 \dots m}$ (table 3)

about splitting



$$a^{(i)} = [\max_{j,j'} (|a_{i,j} - a_{i,j'}|), \max_{j,j'} (|a_{i,j} - a_{i,j'}|), \dots, \max_{j,j'} (|a_{i,j} - a_{i,j'}|)]$$

$$a = \max_i (a_{i1}, a_{i2}, a_{i3}, a_{i4})$$

for each S_{max} .

no statistical significance value.

Learning kernel via Sparse Coding

fundamental point in NL is defining "similarity" function between 2 examples. One way of formalizing this similarity function is by using a kernel function $k: (x^{(1)}, x^{(2)}) \mapsto k(x^{(1)}, x^{(2)})$ that computes a similarity measure.

eg: SVM can be "kernelized" meaning that the algorithm can be written in an equivalent form where it accesses inputs only via computing the kernel function k between two input examples.

eg: SVM "primal"

$$\begin{cases} \text{obj min} & \frac{1}{2} \|w\|^2 + \sum \xi \\ \text{st.} & y_i (w^T x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{cases}$$

$$h(x) = w^T x + b = \sum_{i=1}^n d_i y_i x_i x^T + b$$

Usually, standard kernels are linear, polynomial or gaussian.

SVM "dual"

$$\begin{cases} \text{obj max} & -\frac{1}{2} \sum_i \sum_j d_i d_j y_i y_j x_i^T x_j + \sum d_i \\ \text{st.} & 0 \leq d_i \leq c \\ & \sum d_i y_i = 0 \end{cases}$$

$\langle x_i, x_j \rangle$
 $\uparrow \langle \phi(x_i), \phi(x_j) \rangle$
 ϕ transform function
 \downarrow
 $k(x_i, x_j)$

if we change the kernel in the supervised learning algorithm, we can handle other

definition of similarity:

linear kernel: $k(x^{(r)}, x^{(t)}) = x_{(r)}^T x_{(t)}$

polynomial kernel: $k(x^{(r)}, x^{(t)}) = (x_{(r)}^T x_{(t)} + 1)^d$

gaussian kernel: $k(x^{(r)}, x^{(t)}) = \exp\left(-\frac{\|x^{(r)} - x^{(t)}\|^2}{2\sigma^2}\right)$

RSF

$k(x,x) \geq 0$
symmetric, semi positive
~~semi definite~~
bilinear.

→ Build a specific specialized similarity function (kernel) for the learned representat

idea: instead of comparing input vectors $x^{(r)}$ and $x^{(t)}$, we can compare the generative process for $x^{(r)}$ and $x^{(t)}$: $P(x^{(r)}|b), P(x^{(t)}|b)$.

ie compare the difference in the way the various generative model parameters ^b contribute to generating the input vectors $x^{(r)}$ and $x^{(t)}$.

the method compare the gradient of the log input probabilities with respect to the model parameters

Definition

Fisher score: $U_x = \nabla_b \log(P(x|b))$

Fisher kernel: $k(x^{(r)}, x^{(t)}) = U_{x^{(r)}}^T U_{x^{(t)}}$

$P(x|b) = \int_a P(x|b,a)P(a)da \Rightarrow U_x = \nabla_b \log \int_a P(x|b,a)P(a)da$

point estimate: $\hat{a} = \arg \max_a P(x|b,a)P(a)$ ie find the optimal activations for input x given fixed basis vector b . (MAP).

$U_x = \nabla_b \log(P(x|b, \hat{a})P(\hat{a})) \rightarrow$ ok when derived

Assumption the standard generative model assumes that the reconstruction error

$\eta = x - Ba$ where $B = [b_1, \dots, b_j, \dots, b_d]$

is distributed as a zero-mean Gaussian distribution with covariance $\sigma^2 I$

$P(x|b,a) \propto \exp\left(-\frac{\|\eta\|_2^2}{2\sigma^2}\right)$

let $\hat{\eta} = x - B\hat{a}$

The method compare the gradient of the log input probabilities with respect to the model parameters

Definition

Fisher score : $U_x = \nabla_b \log (P(x|b))$

Fisher kernel : $k(x^{(r)}, x^{(t)}) = U_{x^{(r)}}^T U_{x^{(t)}}$

$P(x|b) = \int_a P(x|b, a) P(a) da \Rightarrow U_x = \nabla_b \log \int_a P(x|b, a) P(a) da$

point estimate : $\hat{a} = \underset{a}{\text{argmax}} P(x|b, a) P(a)$ ie find the optimal activations for input x given fixed basis vector b . (MAP).

...

$U_x = \nabla_b \log (P(x|b, \hat{a}) P(\hat{a}))$ ok when derived

Assumption the standard generative model assumes that the reconstruction error

$\eta = x - Ba$ where $B = [b_1, \dots, b_j, \dots, b_n]$

is distributed as a zero-mean Gaussian distribution with covariance $\sigma^2 I$

$P(x|b, a) \propto \exp(-\frac{\|\eta\|_2^2}{2\sigma^2})$

Let $\hat{\eta} = x - B\hat{a}$

$\log (P(x|b, \hat{a})) \propto \|\hat{\eta}\|_2^2$

So $U_x = \nabla_b \log (P(x|b, \hat{a})) \propto 2\hat{\eta} \nabla_b \hat{\eta} \propto \hat{\eta} \hat{a}$ since \hat{a} is not dependent on b ! constant with respect to a

Finally $U_x \propto \hat{\eta} \hat{a} \Rightarrow U_{x^{(r)}}^T U_{x^{(t)}} \propto \hat{a}^{(r)T} \hat{\eta}^{(r)T} \eta^{(t)} a^{(t)} \propto \hat{a}^{(r)T} a^{(t)} \hat{\eta}^{(r)T} \eta^{(t)}$

$k(x^{(r)}, x^{(t)}) = \underbrace{(\hat{a}^{(r)T} a^{(t)})}_{\text{inner product of the sparse coding features ie linear kernel of sparse coding features}} \underbrace{(\hat{\eta}^{(r)T} \eta^{(t)})}_{\text{inner product of the residuals for input examples}}$ with $\eta = x - \sum_j b_j \hat{a}_j$

optimization

②

$$\text{minimize}_{b, a} \sum_i \|x^{(i)} - \sum_j a_j^{(i)} b_j\|_2^2 + \beta \|a^{(i)}\|_1 \quad \text{st. } \|b_j\|_2^2 \leq 1.$$

convex in B (while holding A fixed), convex in A (while holding B fixed) but not in both simultaneously.

~~not simultaneously~~

⇒ idea: iteratively optimize by alternatingly optimizing with respect to B and A while holding the other fixed.
 (basis) (coefficients)

$$(1) \quad \min_b \|X - BA\|_2^2 \quad B = [b_1 \dots b_j \dots b_n]$$

$$\text{st. } \sum_j b_j^2 \leq 1$$

least squares problem with quadratic constraint

↳ eg: gradient descent with iterative projection.

Here: solve the DUAL.

$$b'_j = b_j - \gamma \nabla F(b) \text{ where } F(b) = \|X - BA\|$$

$$(2) \quad \min_a \|X - BA\|_2^2 + \beta \sum_j |a_j^{(i)}|$$

L1-regularized least square problem.

↳ eg: "feature sign search algorithm".

idea: if we know the signs of the $a_j^{(i)}$ at the optimal value, we can replace each of this term $|a_j^{(i)}|$ with $a_j^{(i)}$ / $-a_j^{(i)}$ or 0. ⇒ unconstrained quadratic optimization problem (QP).

The algorithm is such that it guesses the signs of the coefficients $a_j^{(i)}$ solve a QP and iterate to find optimal \hat{a}_{new} and update the "guess" and iterate.
 or signs accordingly
 ~~before~~
 new