

Introduction to computer vision: Models, Learning, and Inference

Chaloemrat Boonthanawat and Chaiwoot Boonyasiriwat

ABSTRACT

The problem of computer vision can be approached mathematically in term of probabilistic model. The probabilistic context is necessary for the real world image data which are always full with noise and uncertainty. There are three main components of the solution. The first one is the probabilistic model that related the image data to the world state. The model can be classified into two types, discriminative and generative model. The second component is the learning algorithm. The learning algorithm is used to find the parameters in the model using many paired training examples. The third component is the inference algorithm. The inference algorithm is used to infer the world state from the new image data.

INTRODUCTION

The goal of computer vision is to build the machine that can extract useful information from images. The kind of task we are interested in is object classification, image recognition, motion analysis, image restoration, and etc. To put it simply, we seek to automate tasks that human visual system can do.

Mathematically, the problem can be stated as “given the image data \mathbf{x} , the program should be able to infer the world state \mathbf{w} ”. The world state \mathbf{w} may be continuous (the 3D pose of a body model) or discrete (the presence or absence of a particular object). When the world state is continuous, we call regression model and when it is discrete, we call classification model. The image data \mathbf{x} can be discrete (RGB integer value from 0 to 255) or continuous (pixel intensity by real value).

The model of interest is based on probabilistic approach. It is useful because usually the relation between image data and world state is many to one. Moreover, the image data is often full with noise from environment and we need probabilistic model to capture this information.

There are three components for the solution of the problem (Prince, 2012) :

- Model: relates the visual data \mathbf{x} and the world state \mathbf{w} specify by model parameters $\boldsymbol{\theta} = \{\theta_i, \theta_j, \dots\}$.
- Learning algorithm: allows us to fit the parameters $\boldsymbol{\theta}$ using paired training examples $\mathcal{D} = \{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I$.
- Inference algorithm: takes a new observation \mathbf{x}^* and uses the model with a trained parameters $\hat{\boldsymbol{\theta}}$ to return

the posterior probability $\Pr(\mathbf{w}|\mathbf{x}^*, \hat{\boldsymbol{\theta}})$ over the world state \mathbf{w} .

MODELS

Probability distribution and conjugacy

The choices of probability distribution depend on the task of interest. Generally, it is depending on the domain of image data \mathbf{x} or the domain of world state \mathbf{w} (continuous or discrete and the range of possible value) (Table 1.)

Each probability distribution will have parameters ($\boldsymbol{\theta} = \{\theta_1, \theta_2, \dots\}$) corresponding to the model itself. These parameters will be specified by leaning algorithm from paired training examples. In addition, we can also have the probability distribution of parameters from the main model meaning the main model doesn't necessary have a unique parameters. It can have a wide range of possible parameters value. Its parameters will be called hyperparameters.

The choices of probability distribution are chosen such that each is a conjugacy to each other (Table 2.)

The property of conjugacy is that when two distributions multiply, it will give the same distribution with some constant. This property is extremely useful in both learning and inference algorithm by arranging the expression in a closed form.

Table 1: Domain of probability distribution

Distribution	Domain
Bernoulli	$y \in \{0, 1\}$
categorical	$y \in \{1, 2, \dots, K\}$
univariate normal	$y \in \mathbb{R}$
multivariate normal	$\mathbf{y} \in \mathbb{R}^k$

Table 2: Example of possible distribution and its conjugacy

Distribution of model	Distribution of model parameters
Bernoulli	beta
categorical	Dirichlet
univariate normal	normal inverse gamma
multivariate normal	normal inverse Wishart

Hidden variable and more complex probability distribution

Even though, the probability distribution showed in Table 1 is useful for many tasks. It is still limited in real world application because real world image data is complex. To model this type of complex data, we need an additional probability distribution.

The idea of hidden variables is useful for creating more complex distribution from an simpler one. The probability of the model with hidden variables can be written as,

$$\Pr(y) = \int \Pr(y, h)dh \quad (1)$$

$$= \int \Pr(y|h)P(h)dh \quad (2)$$

for continuous case or,

$$\Pr(y) = \sum \Pr(y|h) \quad (3)$$

$$= \sum \Pr(y|h)P(h) \quad (4)$$

for discrete case.

The idea is that we choose probability distribution of $\Pr(y|h)$ and $\Pr(h)$ such that when we compute $\Pr(y)$, it will give the new probability distribution (Table 3.)

Table 3: Complex probability distribution

Distribution	$\Pr(y h)$	$\Pr(h)$
Mixture of Gaussian	multivariate normal	categorical
t -distribution	multivariate normal	gamma
Subspace model	multivariate normal	normal

With the knowledge of probability distribution, every model that related the image data and world state fall into either one of this two categories.

Discriminative model

In this model, we will choose the probability distribution on the world state based on the image data.

$$\Pr(\mathbf{w}|\mathbf{x}, \boldsymbol{\theta}_m) \quad (5)$$

The advantage of this approach is that the inference algorithm is extremely easy by just using the new image data x^* in the discriminative model.

Generative model

In this model, we will choose the probability distribution on the image data based on the world state. We have to model both the probability distribution between image data and world state and also the probability distribution

of the world state alone.

$$\Pr(\mathbf{x}|\mathbf{w}, \boldsymbol{\theta}_m) \quad (6)$$

$$\Pr(\mathbf{w}|\boldsymbol{\theta}_w) \quad (7)$$

The advantage of this approach is that we can incorporate the knowledge about how world state generate the image data into the model. But the disadvantage is that inference algorithm is more complicated for computing.

LEARNING

There are three main ways we can used to fit the parameters $\boldsymbol{\theta}$ by using I paired training examples $\mathcal{D} = \{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I$. The assumption we used is that each pair of data and world state is independent from each other.

Maximum likelihood

The maximum likelihood (ML) method finds the set of parameters $\hat{\boldsymbol{\theta}}$ in which paired training examples $\mathcal{D} = \{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I$ are most likely. For a discriminative model, we can write the equation as,

$$\hat{\boldsymbol{\theta}}_m = \operatorname{argmax}_{\boldsymbol{\theta}_m} [\Pr(\mathbf{w}_1, \dots, \mathbf{w}_I | \mathbf{x}_1, \dots, \mathbf{x}_I, \boldsymbol{\theta}_m)] \quad (8)$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}_m} \left[\prod_{i=1}^I \Pr(\mathbf{w}_i | \mathbf{x}_i, \boldsymbol{\theta}_m) \right] \quad (9)$$

For a generative model, we can write the equation as,

$$\hat{\boldsymbol{\theta}}_m = \operatorname{argmax}_{\boldsymbol{\theta}_m} [\Pr(\mathbf{x}_1, \dots, \mathbf{x}_I | \mathbf{w}_1, \dots, \mathbf{w}_I, \boldsymbol{\theta}_m)] \quad (10)$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}_m} \left[\prod_{i=1}^I \Pr(\mathbf{x}_i | \mathbf{w}_i, \boldsymbol{\theta}_m) \right] \quad (11)$$

and for parameters of the world state probability distribution $\Pr(\mathbf{w}|\boldsymbol{\theta}_w)$, we use the training world state $\{\mathbf{w}_i\}_{i=1}^I$

$$\hat{\boldsymbol{\theta}}_w = \operatorname{argmax}_{\boldsymbol{\theta}_w} [\Pr(\mathbf{w}_1, \dots, \mathbf{w}_I | \boldsymbol{\theta}_w)] \quad (12)$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}_w} \left[\prod_{i=1}^I \Pr(\mathbf{w}_i | \boldsymbol{\theta}_w) \right] \quad (13)$$

Maximum a posteriori

In maximum a posteriori (MAP), we introduce prior information about the parameters $\boldsymbol{\theta}$. These may come from previous experience or knowledge. For a discriminative

model, we can write the equation as,

$$\hat{\theta}_m = \operatorname{argmax}_{\theta_m} [\Pr(\theta_m | \{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I)] \quad (14)$$

$$= \operatorname{argmax}_{\theta_m} [\Pr(\{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I | \theta_m) \Pr(\theta_m)] \quad (15)$$

$$= \operatorname{argmax}_{\theta_m} \left[\prod_{i=1}^I \Pr(\mathbf{w}_i, \mathbf{x}_i | \theta_m) \Pr(\theta_m) \right] \quad (16)$$

$$= \operatorname{argmax}_{\theta_m} \left[\prod_{i=1}^I \Pr(\mathbf{w}_i | \mathbf{x}_i, \theta_m) \Pr(\mathbf{x}_i | \theta_m) \Pr(\theta_m) \right] \quad (17)$$

$$= \operatorname{argmax}_{\theta_m} \left[\prod_{i=1}^I \Pr(\mathbf{w}_i | \mathbf{x}_i, \theta_m) \Pr(\theta_m) \right] \quad (18)$$

For a generative model, we can write the equation as,

$$\hat{\theta}_m = \operatorname{argmax}_{\theta_m} [\Pr(\theta_m | \{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I)] \quad (19)$$

$$= \operatorname{argmax}_{\theta_m} \left[\prod_{i=1}^I \Pr(\mathbf{w}_i, \mathbf{x}_i | \theta_m) \Pr(\theta_m) \right] \quad (20)$$

$$= \operatorname{argmax}_{\theta_m} \left[\prod_{i=1}^I \Pr(\mathbf{x}_i | \mathbf{w}_i, \theta_m) \Pr(\mathbf{w}_i | \theta_m) \Pr(\theta_m) \right] \quad (21)$$

$$= \operatorname{argmax}_{\theta_m} \left[\prod_{i=1}^I \Pr(\mathbf{x}_i | \mathbf{w}_i, \theta_m) \Pr(\theta_m) \right] \quad (22)$$

We can see that maximum likelihood is a special case of maximum a posteriori when we don't have any information about parameters in advance.

In order to simplify the learning algorithm, we can use the logarithm of the expression because it is monotonic function. The maximum parameters of transformed function will still remain the same, but it will give an easier expression to deal with.

Expectation maximization algorithm

The EM algorithm is used to find the parameters $\hat{\theta}_m$ with the equation of the form (maximum likelihood),

$$\hat{\theta}_m = \operatorname{argmax}_{\theta_m} \left\{ \sum_{i=1}^I \log \left[\sum \Pr(\mathbf{x}_i, \mathbf{w}_i | \mathbf{h}_i, \theta_m) \Pr(\mathbf{h}_i) \right] \right\} \quad (23)$$

or in this form (maximum a posteriori)

$$= \operatorname{argmax}_{\theta_m} \left\{ \sum_{i=1}^I \log \left[\sum \Pr(\mathbf{x}_i, \mathbf{w}_i | \mathbf{h}_i, \theta_m) \Pr(\mathbf{h}_i) \right] \Pr(\theta_m) \right\} \quad (24)$$

The same is true for the case when hidden variables are continuous. There are two main steps in the algorithm which will be iterated until the maximum parameters are found. 1.) the E-step

$$q_i(\mathbf{h}_j) = \frac{\Pr(\mathbf{x}_i, \mathbf{w}_i | \mathbf{h}_j, \theta_m^{[t]}) \Pr(\mathbf{h}_j | \theta_m^{[t]})}{\Pr(\mathbf{x}_i, \mathbf{w}_i)} \quad (25)$$

2.) the M-step

$$\theta_m^{[t+1]} = \operatorname{argmax}_{\theta_m} \left\{ \sum_{i=1}^I \sum q_i(\mathbf{h}_j) \log [\Pr(\mathbf{x}_i, \mathbf{w}_i, \mathbf{h}_i | \theta_m)] \right\} \quad (26)$$

The Bayesian approach

The Bayesian approach is similar to the maximum a posteriori. Instead of finding a unique parameters $\hat{\theta}$, we find the probability of all parameters. Each parameters will contribute to the inference based on its probability, the expression is

$$\Pr(\theta | \mathcal{D}) = \Pr(\theta | \{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I) \quad (27)$$

$$= \frac{\Pr(\{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I | \theta) \Pr(\theta)}{\Pr(\{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I)} \quad (28)$$

$$= \frac{\prod_{i=1}^I \Pr(\mathbf{w}_i, \mathbf{x}_i | \theta) \Pr(\theta)}{\Pr(\{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I)} \quad (29)$$

INFERENCE

In the maximum likelihood (ML) and maximum a posteriori method (MAP), we would get a unique parameters $\hat{\theta}$. These parameters is used to find the probability of each possible world state \mathbf{w} given the new image data \mathbf{x}^* .

Inference algorithm for discriminative model (ML, MAP method)

For discriminative model, we have already directly constructed an expression for the posterior distribution, and we simply evaluate it with the new data.

$$\Pr(\mathbf{w} | \mathbf{x}^*, \hat{\theta}_m) \quad (30)$$

Inference algorithm for generative model (ML, MAP method)

For generative model, we use Bayes' rule to calculate the posterior distribution $\Pr(\mathbf{w} | \mathbf{x}^*, \hat{\theta})$. For continuous case, we got

$$\Pr(\mathbf{w} | \mathbf{x}^*, \hat{\theta}) = \frac{\Pr(\mathbf{x}^* | \mathbf{w}, \hat{\theta}_m) \Pr(\mathbf{w} | \hat{\theta}_m)}{\int \Pr(\mathbf{x}^* | \mathbf{w}, \hat{\theta}_m) \Pr(\mathbf{w} | \hat{\theta}_m) d\mathbf{w}} \quad (31)$$

For discrete case, we got

$$\Pr(\mathbf{w} | \mathbf{x}^*, \hat{\theta}) = \frac{\Pr(\mathbf{x}^* | \mathbf{w}, \hat{\theta}_m) \Pr(\mathbf{w} | \hat{\theta}_m)}{\sum_{\mathbf{w}} \Pr(\mathbf{x}^* | \mathbf{w}, \hat{\theta}_m) \Pr(\mathbf{w} | \hat{\theta}_m)} \quad (32)$$

Inference algorithm for Bayesian approach

Evaluating the predictive distribution is more difficult for the Bayesian case since we have not estimated a unique parameter for a model. But we instead have a probability distribution over all possible models (all possible parameters). To calculate the probability of world state given the new data \mathbf{x}^* , we need to weight the probability of all possible models. For a discriminative model, we can write the equation as (Bishop, 2006),

$$\begin{aligned} \Pr(\mathbf{w}^*|\mathbf{x}^*) &= \int \Pr(\boldsymbol{\theta}_m|\mathcal{D}) \Pr(\mathbf{w}^*|\mathbf{x}^*, \boldsymbol{\theta}_m) d\boldsymbol{\theta}_m \\ &= \int \left[\frac{\prod_{i=1}^I \Pr(\mathbf{w}_i, \mathbf{x}_i|\boldsymbol{\theta}_m) \Pr(\boldsymbol{\theta}_m)}{\Pr(\{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I)} \right] \Pr(\mathbf{w}^*|\mathbf{x}^*, \boldsymbol{\theta}_m) d\boldsymbol{\theta}_m \\ &= \int \left[\frac{\prod_{i=1}^I \Pr(\mathbf{w}_i|\mathbf{x}_i, \boldsymbol{\theta}_m) \Pr(\boldsymbol{\theta}_m)}{\prod_{i=1}^I \Pr(\mathbf{w}_i|\mathbf{x}_i)} \right] \Pr(\mathbf{w}^*|\mathbf{x}^*, \boldsymbol{\theta}_m) d\boldsymbol{\theta}_m \end{aligned} \quad (33)$$

For a generative model, we can write the equation as,

$$\Pr(\mathbf{w}^*|\mathbf{x}^*) = \int \Pr(\boldsymbol{\theta}|\mathcal{D}) \Pr(\mathbf{w}^*|\mathbf{x}^*, \boldsymbol{\theta}) d\boldsymbol{\theta} \quad (34)$$

$$\begin{aligned} &= \int \Pr(\boldsymbol{\theta}_m, \boldsymbol{\theta}_w|\{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I) \\ &\quad \left[\frac{\Pr(\mathbf{x}^*|\mathbf{w}, \boldsymbol{\theta}_m) \Pr(\mathbf{w}|\boldsymbol{\theta}_w)}{\int \Pr(\mathbf{x}^*|\mathbf{w}, \boldsymbol{\theta}_m) \Pr(\mathbf{w}|\boldsymbol{\theta}_w) d\mathbf{w}} \right] d\boldsymbol{\theta}_m d\boldsymbol{\theta}_w \end{aligned} \quad (35)$$

$$\begin{aligned} &= \int \Pr(\boldsymbol{\theta}_m|\{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I) \Pr(\boldsymbol{\theta}_w|\{\mathbf{x}_i, \mathbf{w}_i\}_{i=1}^I) \\ &\quad \left[\frac{\Pr(\mathbf{x}^*|\mathbf{w}, \boldsymbol{\theta}_m) \Pr(\mathbf{w}|\boldsymbol{\theta}_w)}{\int \Pr(\mathbf{x}^*|\mathbf{w}, \boldsymbol{\theta}_m) \Pr(\mathbf{w}|\boldsymbol{\theta}_w) d\mathbf{w}} \right] d\boldsymbol{\theta}_m d\boldsymbol{\theta}_w \end{aligned} \quad (36)$$

$$\begin{aligned} &= \int \left[\frac{\prod_{i=1}^I \Pr(\mathbf{x}_i|\mathbf{w}_i, \boldsymbol{\theta}_m) \Pr(\boldsymbol{\theta}_m)}{\prod_{i=1}^I \Pr(\mathbf{x}_i|\mathbf{w}_i)} \right] \left[\prod_{i=1}^I \Pr(\mathbf{w}_i|\boldsymbol{\theta}_w) \Pr(\boldsymbol{\theta}_w) \right] \\ &\quad \left[\frac{\Pr(\mathbf{x}^*|\mathbf{w}, \boldsymbol{\theta}_m) \Pr(\mathbf{w}|\boldsymbol{\theta}_w)}{\int \Pr(\mathbf{x}^*|\mathbf{w}, \boldsymbol{\theta}_m) \Pr(\mathbf{w}|\boldsymbol{\theta}_w) d\mathbf{w}} \right] d\boldsymbol{\theta}_m d\boldsymbol{\theta}_w \end{aligned} \quad (37)$$

In practice, if we choose an appropriate probability distribution (conjugacy), we do not necessary have to calculate the integral directly.

CONCLUSION

Even though the three components are the foundation of the solution, there are still many issues about performance, robustness and the amount of training data required for learning algorithm. More technique and method is needed to ensure that the program work nicely in real life. Nevertheless, these knowledge provided a foundation that can be applied to any task at hand.

ACKNOWLEDGMENTS

We would like to thank members of the Mahidol University Center for Scientific Computing (MCSC) group for comments and suggestions.

REFERENCES

- Bishop, C. M., 2006, Pattern recognition and machine learning: Springer.
 Prince, S. J. D., 2012, Computer vision: Models, learning, and inference: Cambridge University Press, New York.